

---

# Bandit Multi-linear DR-Submodular Maximization and Its Applications on Adversarial Submodular Bandits

---

Zongqi Wan<sup>1,2</sup> Jialin Zhang<sup>1,2</sup> Wei Chen<sup>3</sup> Xiaoming Sun<sup>1,2</sup> Zhijie Zhang<sup>4</sup>

## Abstract

We investigate the online bandit learning of the monotone multi-linear DR-submodular functions, designing the algorithm `BanditMLSM` that attains  $O(T^{2/3} \log T)$  of  $(1 - 1/e)$ -regret. Then we reduce submodular bandit with partition matroid constraint and bandit sequential monotone maximization to the online bandit learning of the monotone multi-linear DR-submodular functions, attaining  $O(T^{2/3} \log T)$  of  $(1 - 1/e)$ -regret in both problems, which improve the existing results. To the best of our knowledge, we are the first to give a sublinear regret algorithm for the submodular bandit with partition matroid constraint. A special case of this problem is studied by Streeter et al. (2009). They prove a  $O(T^{4/5})$   $(1 - 1/e)$ -regret upper bound. For the bandit sequential submodular maximization, the existing work proves an  $O(T^{2/3})$  regret with a suboptimal  $1/2$  approximation ratio (Niazadeh et al., 2021).

## 1. Introduction

Research on multi-armed bandit problems has developed rapidly in the last two decades. After the classical finite-arm bandit and linear bandit were well-studied in both stochastic and adversarial settings, people were starting to consider more general bandit problems. Submodular bandit is such an object being considered due to its ability to characterize the diminishing property of reward function in realistic applications.

In submodular bandit, an optimizer/decision-maker needs to select a feasible subset each round, then a monotone submodular reward function of this round is determined

---

<sup>1</sup>Institute of Computing Technology, Chinese Academy of Sciences <sup>2</sup>University of Chinese Academy of Sciences <sup>3</sup>Microsoft Research <sup>4</sup>Center for Applied Mathematics of Fujian Province, School of Mathematics and Statistics, Fuzhou University. Correspondence to: Zhijie Zhang <zhang@fzu.edu.cn>.

stochastically or adversarially, and the decision maker obtains a reward according to the reward function of this round. In this paper, we consider the adversarial setting. In the adversarial submodular bandit literature, the meta-action technique proposed by Streeter & Golovin (2008) is commonly used to obtain a sublinear regret algorithm. This technique employs online optimizers for each offline step of an offline algorithm to mimic it in an online manner. This technique reaches  $O(T^{2/3})$   $(1 - 1/e)$ -regret with the cardinality constraint (Streeter & Golovin, 2008). Subsequently, Streeter et al. applied this technique to the assignment constraint (Streeter et al., 2009), which is a special partition matroid where the feasible set can only select one item from each partition. They obtained an  $O(T^{4/5})$   $(1 - 1/e)$ -regret in this situation. A recent work (Niazadeh et al., 2021) uses a Blackwell algorithm to turn offline greedy algorithms into online regret minimization algorithms. As an application, they reproduced the  $O(T^{2/3})$   $(1 - 1/e)$ -regret for the cardinality constraint.

In this paper, we present a different approach for the adversarial submodular bandit. We reduce the submodular bandit into a bandit multi-linear DR-submodular maximization problem. The DR-submodular function is a kind of non-convex function with theoretical guarantees in optimization, which has received much attention in recent years (Bian et al., 2017a;b; Niazadeh et al., 2020). There are several works considering the online full information or bandit feedback learning of the DR-submodular function (Chen et al., 2018; Zhang et al., 2019; Raut et al., 2020; Thang & Srivastav, 2021; Zhang et al., 2022a;b). DR-submodularity is inspired by the submodular set function, and we find it useful for designing submodular bandit algorithms due to its continuity. Specifically, we propose the function class called the multi-linear DR-submodular function. A multi-linear DR-submodular function is a DR-submodular function, and we additionally require it to be a multi-variable polynomial with the degree of each variable not exceeding 1. We propose the algorithm `BanditMLSM` for the bandit maximization of this function class and reach the  $(1 - 1/e)$ -regret of  $\tilde{O}(T^{2/3})$ , which is far better than the  $O(T^{5/6})$   $(1 - 1/e)$ -regret bound achieved on the general bandit DR-submodular maximization problem (Niazadeh et al., 2021). Multi-linear DR-submodular function captures the property of the multi-

linear extension of a submodular set function. In fact, a multi-linear extension is a special case of multi-linear DR-submodular functions.

Our next goal is to reduce discrete submodular bandit to bandit multi-linear DR-submodular maximization problem. A natural idea is to run `BanditMLSM` on the multi-linear extension of a submodular set function. However, this idea fails to reduce the submodular bandit to bandit multi-linear DR-submodular maximization problem. This is because the function value of the multi-linear extension cannot be estimated unbiasedly while the constraint is not trivial, which is because the value of the multi-linear extension may require obtaining feedback on a set function value  $f(S)$  where the set  $S$  is outside the constraint (e.g. the cardinality constraint). This is not allowed in the bandit feedback model. To address this issue, we propose a new kind of continuous extension which is also multi-linear DR-submodular. Then we run `BanditMLSM` on that extension. We try our continuous approach on submodular bandit with partition matroid constraint and bandit sequential submodular maximization, generalizing and improving the previous results, see [Section 1.2](#).

**More related works** There is also some research on the stochastic submodular bandit. A model named linear submodular bandit has been studied ([Yue & Guestrin, 2011](#); [Chen et al., 2017](#)). The model assumes that the reward function is a linear combination of several known submodular functions, only the weights of each submodular function are unknown to the decision maker, and the model requires the noisy marginal gain as the stochastic feedback. Many studies focus on the online influence maximization problem ([Vaswani et al., 2015](#); [Chen et al., 2016](#); [Wang & Chen, 2017](#); [Wu et al., 2019](#); [Li et al., 2020](#); [Zhang et al., 2022c](#)), where the submodular reward function is induced by an information diffusion process on a social network. In this problem, different feedback models are studied. However, all the studies above assume extra information more than a *full-bandit* feedback model where the decision maker can only observe the reward of the action played. We only notice two works studying the full-bandit feedback model: ([Nie et al., 2022](#); [2023](#)). [Nie et al. \(2022\)](#) studied the bandit monotone submodular maximization with cardinality constraint, attaining  $(1 - 1/e)$ -regret of order  $O(T^{2/3})$ ; [Nie et al. \(2023\)](#) design a framework which adapts an  $\alpha$ -approximate offline algorithm into a stochastic bandit algorithm with  $O(T^{2/3}(\log(T))^{1/3})$   $\alpha$ -regret. The framework needs the offline algorithm to be robust to small errors. Besides the above, [Foster & Rakhlin \(2021\)](#) studied the submodular contextual bandit.

**Remark on the stochastic submodular bandit** While [Nie et al. \(2022\)](#) make in their paper a weaker assumption

that the online reward function need not be a monotone submodular but only need to be monotone submodular in expectation, we find our algorithm can also be applied in this setting even our adversarial submodular bandit model requires the reward function to be monotone submodular. The key observation is, when we apply the adversarial submodular bandit problem on a stochastic submodular bandit environment, we should see the expected submodular function rather than the stochastically realized reward function as the online reward function selected by the adversary, and see the stochastic feedback as an unbiased estimate of the true value of the expected function. We will explain this in [Appendix F](#).

### 1.1. Bandit Optimization Model

Adversarial bandit optimization problems can be formalized as a repeated game between an optimizer and an adversary. The game lasts for  $T$  rounds and  $T$  is known to both players. In  $t$ -th round, the optimizer chooses an action  $x_t$  from an action set  $\mathcal{K}$ , then the adversary chooses a reward function  $f_t \in \mathcal{F}$ . The action set  $\mathcal{K}$  and the reward function set  $\mathcal{F}$  are determined by specific bandit problems. Generally,  $f_t$  maps  $\mathcal{K}$  to a bounded interval  $[0, M] \subseteq \mathbb{R}$ . The optimizer gets reward  $f_t(x_t)$  and it can only observe the value  $f_t(x_t)$ , which is called the *bandit* feedback model. Sometimes people also call it the *full-bandit* model to distinguish it from the semi-bandit model where the optimizer can observe more information, in this paper bandit and full-bandit are the same thing.

In this paper, we consider *oblivious* adversary, which means the reward functions  $f_t$  can not be adaptively selected according to  $x_1, x_2, \dots, x_t$ . In other words, we can think the adversary selects  $f_t \in \mathcal{F}$  for each  $1 \leq t \leq T$  before the game starts and these functions are not revealed to the optimizer. Our goal is to design a strategy for the optimizer to minimize its cumulative  $\alpha$ -regret during  $T$  rounds,

$$\mathcal{R}_\alpha(T) = \max_{x^* \in \mathcal{K}} \mathbb{E} \left[ \sum_{t=1}^T (\alpha f_t(x^*) - f_t(x_t)) \right].$$

The action set  $\mathcal{K}$  could be some structured set, maybe infinite or finite size. For the convenience of subsequent descriptions, we use  $\mathcal{S}$  to denote the finite action set of the optimizer and use  $\mathcal{K}$  to denote the infinite action set. Given  $\mathcal{K}$  and  $\mathcal{F}$ , we call the game a  $(\mathcal{K}, \mathcal{F})$ -bandit.

When we consider the bandit multi-linear monotone DR-submodular maximization and bandit DR-submodular maximization in [Section 3](#) and [Section 4](#), we focus on the situation that  $\mathcal{K}$  satisfies [Assumption 1.1](#).

**Assumption 1.1.** We assume  $\mathcal{K}$  is a compact convex subset of  $\mathbb{R}^d$  containing  $\mathbf{0}$ , and  $\mathcal{K} \subseteq D\mathbb{B}_d$  for some constant  $D$ ,

where  $\mathbb{B}_d$  is a  $d$ -dimensional unit ball.

Before further describing the model we are considering, we give several definitions.

**Definition 1.2** (Monotonicity). There is a natural partial order on  $\mathbb{R}^d$ . For  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ , if  $x_i \geq y_i \forall i \in [d]$ , then  $\mathbf{x} \geq \mathbf{y}$ . For a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , if for any  $\mathbf{x} \geq \mathbf{y}$ ,  $f(\mathbf{x}) \geq f(\mathbf{y})$ , we call  $f$  a monotone function.

**Definition 1.3** (DR-submodularity). Let  $\mathcal{X} = \prod_{i=1}^d \mathcal{X}_i$  be a subset of  $\mathbb{R}^d$ , where  $\mathcal{X}_i$  is an interval  $[0, a_i]$ . A continuous function  $f : \mathcal{X} \rightarrow \mathbb{R}_+$  is called a DR-submodular function if for any  $\mathbf{x} \geq \mathbf{y}$ ,  $\lambda \in \mathbb{R}^+$ , and the  $i$ -th base vector  $\mathbf{e}_i$  for any  $i \in [d]$ ,

$$f(\mathbf{x} + \lambda \mathbf{e}_i) - f(\mathbf{x}) \leq f(\mathbf{y} + \lambda \mathbf{e}_i) - f(\mathbf{y}).$$

Moreover, if  $f$  is second-order differentiable, then the DR-submodularity is equivalent to  $\frac{\partial^2 f}{\partial x_i \partial x_j} \leq 0, \forall i, j \in [d]$ .

**Definition 1.4** (Multi-linearity). We say function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  a multi-linear function if  $f$  is polynomial of  $d$  variables, and for any variable  $x_i$ , the degree of  $x_i$  in each term of  $f$  is no more than 1.

**Definition 1.5** (Lipschitz condition and smoothness). Let  $\|\cdot\|$  be the  $L^2$ -norm. For continuous differentiable function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , if  $|f(\mathbf{x}) - f(\mathbf{y})| \leq L_1 \|\mathbf{x} - \mathbf{y}\|$  for any  $\mathbf{x}, \mathbf{y}$ , we say  $f$  is  $L_1$ -lipschitz continuous. If  $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L_2 \|\mathbf{x} - \mathbf{y}\|$ , we say  $f$  is  $L_2$ -smooth.

With the above definitions, we consider two reward function sets in Section 3 and Section 4:

- $\mathcal{F}_{DS}$ : The set of monotone DR-submodular functions, which are  $L_1$ -lipschitz continuous,  $L_2$ -smooth, and  $f(\mathbf{0}) = 0$ .
- $\mathcal{F}_{MDS}$ : The set of monotone multi-linear DR-submodular functions, which are  $L_1$ -lipschitz continuous, and  $f(\mathbf{0}) = 0$ .

## 1.2. Our Results

We are the first to consider the  $(\mathcal{K}, \mathcal{F}_{MDS})$ -bandit, i.e. Bandit Monotone Multi-linear DR-Submodular Maximization (BMMDSM). We observe that the gradient of multi-linear functions can be written as a linear combination of finite function values. Therefore, compared to the standard one-point gradient estimator proposed in (Flaxman et al., 2005) which is used in previous works (Zhang et al., 2019; Niazadeh et al., 2021), we propose a better one-point gradient estimator for monotone multi-linear DR-submodular functions. Along with other techniques including self-concordant barrier and non-oblivious technique, we propose the algorithm `BanditMLSM`, which achieves

$(1 - 1/e)$ -regret of  $\tilde{O}(T^{2/3})$ . Here the  $\tilde{O}$  hides the  $\log T$  factor.

As a secondary result, we also improved the  $(1 - 1/e)$ -regret of general  $(\mathcal{K}, \mathcal{F}_{DS})$ -bandit. This bandit is studied in (Zhang et al., 2019; Niazadeh et al., 2021), where they gave the  $(1 - 1/e)$ -regret bounds of  $O(T^{8/9})$  and  $O(T^{5/6})$  respectively. We proposed the algorithm `BanditDRSM` which achieves the  $(1 - 1/e)$ -regret of  $\tilde{O}(T^{3/4})$ . Compared with their assumptions on functions, we add a new assumption that  $f_t(\mathbf{0}) = 0$ . Fortunately, this assumption is satisfied by many applications of DR-submodular maximization, including *optimal budget allocation with continuous assignment*, *senser placement*, *softmax extension* and so on (Bian et al., 2017a;b). For the constraint set, they assume  $\mathcal{K}$  is downward closed while we do not make this assumption.

Our main contribution is to propose a continuous approach for combinatorial full-bandit, for example, the case where online functions are submodular set functions and the constraint is a partition matroid. Talking about the continuous approach, a natural idea is reducing combinatorial bandit to Bandit Monotone Multi-linear DR-submodular Maximization using the classical multi-linear extension technique. For a submodular set function  $g$  over the ground set  $G = \{1, 2, \dots, n\}$ , its multi-linear extension  $f : [0, 1]^n \rightarrow \mathbb{R}^+$  is defined as

$$f(\mathbf{x}) = \sum_{S \subseteq G} g(S) \prod_{i \in S} x_i \prod_{i \notin S} (1 - x_i).$$

From the definition of the multi-linear extension, we can see that it needs the value information of set function  $g$  over all subsets of  $G$ . However, in the submodular bandit, one can only take the action which satisfies the constraint, thus the algorithm can not explore the value information outside the constraint. As a result, Zhang et al. (2019) proved that it is impossible to construct an unbiased estimate of  $f$  and the gradient of  $f$ . That is to say, classical multi-linear extension is not a good candidate for our goal.

To overcome the above difficulties, we propose other continuous multi-linear DR-submodular extensions which require only the information of the feasible action. We select two submodular bandit problems to clarify our methodology: Bandit Monotone Submodular Maximization with Partition Matroid Constraint (BMSMPM) and Bandit Sequential Submodular Maximization (BSSM). The results are summarized in Table 1. Previous works have studied two special cases of BMSMPM: cardinality constraint (Streeter & Golovin, 2008) and the assignment problem (Streeter et al., 2009). We improve the regret bound of the bandit assignment problem and reproduce an  $\tilde{O}(T^{2/3})$   $(1 - 1/e)$ -regret for cardinality constraint. To our best knowledge, we are the first to give a sublinear  $(1 - 1/e)$ -regret algorithm for bandit monotone submodular maximization with general partition matroid

Table 1. Our results comparing to the previous results.

| Problem  | Our $\alpha$ and regret  | Previous $\alpha$ and regret   |
|--|--|--|
| Bandit DR-submodular maximization results                      |  |  |
| Bandit Multi-linear Monotone DR-Submodular Maximization        | Theorem 3.3<br>$1 - 1/e, O(d^{4/3}T^{2/3} \log(T))$  | \  |
| Bandit Monotone DR-Submodular Maximization                     | Theorem 4.1<br>$1 - 1/e, O(d^{1/2}T^{3/4} \log(T))$  | (Niazadeh et al., 2021)<br>$1 - 1/e, O(d(\log(d))^{1/6}T^{5/6})$       |
| Applications on adversarial submodular bandits                 |  |  |
| Bandit Assignment Problem                                      | Corollary 5.4 <sup>†</sup><br>$1 - 1/e, O( G ^{5/3}T^{2/3} \log(T))$                               | (Streeter et al., 2009)<br>$1 - 1/e, O(T^{4/5})$ <sup>‡</sup>          |
| Bandit Monotone Submodular Maximization over Partition Matroid | Corollary 5.4<br>$1 - 1/e, O\left(\left(\sum_{k=1}^K r_k  G_k \right)^{5/3} T^{2/3} \log T\right)$ | \  |
| Bandit Sequential Submodular Maximization <sup>¶</sup>         | Corollary 5.6<br>$1 - 1/e, O( G ^{10/3}T^{2/3} \log(T))$   | (Niazadeh et al., 2021)<br>$1/2, O( G ^{5/3}(\log( G ))^{1/3}T^{2/3})$ |

<sup>†</sup> Bandit assignment problem is a special case of Bandit Monotone Submodular Maximization over Partition Matroid where  $r_k = 1$  and  $G_k = G$ , so this regret bound can be directly derived from Corollary 5.4. <sup>‡</sup> In the original paper (Streeter et al., 2009), the regret is written in the form which contains the optimal cumulative reward value, which will continue to be bounded to  $T$  usually, leading to a bad dependent on  $T$ . So we re-trade off their regret and write it in terms of  $T$  so that it can be compared with our regret bound. <sup>¶</sup> Compared with the setting in (Niazadeh et al., 2021), we actually add a new assumption that there is a dummy element in the ground set that always has 0 marginal gain. This assumption can be satisfied easily in realistic applications, see Section 5.3.

constraint. BSSM is motivated by maximizing user engagement on online retailing platforms. It is first studied in (Niazadeh et al., 2021), and their algorithm attains  $O(T^{2/3})$   $1/2$ -regret while the  $1/2$  approximation ratio is not tight. We improve this result to  $\tilde{O}(T^{2/3})$   $(1 - 1/e)$ -regret, leading to the tight approximation ratio.

In summary, we make the following contributions:

- We are the first to study the bandit maximization of multi-linear monotone DR-submodular functions and propose a  $\tilde{O}(T^{2/3})$   $(1 - 1/e)$ -regret algorithm.
- We improve the previous result of bandit maximization of general monotone DR-submodular functions to  $\tilde{O}(T^{3/4})$   $(1 - 1/e)$ -regret by better exploiting the smoothness.
- We propose a continuous approach to reducing combinatorial bandit to multi-linear DR-submodular bandit. Using this continuous approach, we propose the first sublinear regret algorithm for submodular bandit with partition matroid constraint, which also improves the result of a previous work (Streeter et al., 2009) that studied the special case of this problem. We also improve the previous approximation ratio of Bandit Sequential Submodular Maximization from  $1/2$  to tight  $1 - 1/e$ .

## 2. Preliminary

### 2.1. Regularized Follow the Leader and Self-Concordant Functions

Regularized Follow The Leader (RFTL) is a commonly used algorithm for online optimization. While applying on a sequence of vector  $\{\mathbf{g}_q\}_{q=1}^Q$  with constraint  $\mathcal{K}$ , RFTL outputs a sequence of point  $\{\mathbf{x}_q\}_{q=1}^Q$ , where

$$\begin{aligned} \mathbf{x}_1 &= \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}} \Phi(\mathbf{x}) \\ \mathbf{x}_{q+1} &= \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}} \left( \eta \sum_{s=1}^q \langle -\mathbf{g}_s, \mathbf{x} \rangle + \Phi(\mathbf{x}) \right). \end{aligned}$$

Here  $\Phi(\mathbf{x})$  is an arbitrary regularizer,  $\eta$  is a parameter. In this paper, we use a self-concordant barrier of  $\mathcal{K}$  as the regularizer of RFTL. Self-concordant barrier was first proposed in convex optimization literature, and it was introduced to the bandit optimization problem in (Abernethy et al., 2008).

**Definition 2.1** (Self-concordant Barrier (Hazan et al., 2016)). Let  $\mathcal{K} \in \mathbb{R}^d$  be a convex set with non empty interior  $\operatorname{int}(\mathcal{K})$ . We call the function  $\Phi : \operatorname{int}(\mathcal{K}) \rightarrow \mathbb{R}$  a  $\nu$ -self-concordant barrier of  $\mathcal{K}$  if:

- (1)  $\Phi$  is three-times continuously differentiable, convex, and approaches infinity along any sequence of points approaching the boundary of  $\mathcal{K}$ ;
- (2) For every  $\mathbf{h} \in \mathbb{R}^d$  and  $\mathbf{x} \in \operatorname{int}(\mathcal{K})$ , the following holds:  $|\nabla^3 \Phi(\mathbf{x})[\mathbf{h}, \mathbf{h}, \mathbf{h}]| \leq 2(\nabla^2 \Phi(\mathbf{x})[\mathbf{h}, \mathbf{h}])^{3/2}$ ,

$|\nabla\Phi(x)[\mathbf{h}]| \leq \nu^{1/2}(\nabla^2\Phi(\mathbf{x})[\mathbf{h}, \mathbf{h}])^{1/2}$ . where the third-order differential is defined as  $\nabla^3\Phi(\mathbf{x})[\mathbf{h}, \mathbf{h}, \mathbf{h}] := \frac{\partial^3}{\partial t_1 \partial t_2 \partial t_3} \Phi(x + t_1 \mathbf{h} + t_2 \mathbf{h} + t_3 \mathbf{h})|_{t_1=t_2=t_3=0}$ .

**Definition 2.2** (Local norm). The Hessian of self-concordant barrier induces a local norm at every  $x \in \text{int}(\mathcal{K})$ , denoted as  $\|\cdot\|_{\Phi, \mathbf{x}}$ . We denote its dual norm as  $\|\cdot\|_{\Phi, \mathbf{x}, *}$ . For any  $\mathbf{v} \in \mathbb{R}^d$ ,

$$\begin{aligned} \|\mathbf{v}\|_{\Phi, \mathbf{x}} &= \sqrt{\mathbf{v}^T \nabla^2 \Phi(\mathbf{x}) \mathbf{v}} \\ \|\mathbf{v}\|_{\Phi, \mathbf{x}, *} &= \sqrt{\mathbf{v}^T (\nabla^2 \Phi(\mathbf{x}))^{-1} \mathbf{v}}. \end{aligned}$$

The following theorem is proved in (Abernethy et al., 2008). It shows that, if we set the regularizer to be a self-concordant barrier of  $\mathcal{K}$  and the algorithm can access the unbiased estimator of  $g_t$ , then the regret of the generated solution sequence  $\{x_q\}_{q=1}^Q$  can be bounded in terms of the local norm of the estimator.

**Theorem 2.3** ((Abernethy et al., 2008)). *Let  $\mathcal{K}$  be a convex set,  $\Phi(x)$  be a self-concordant barrier on  $\mathcal{K}$ ,  $\{\tilde{\mathbf{g}}_q\}_{q=1}^Q$  be a vector sequence. If  $\tilde{\mathbf{g}}_q$  is an unbiased estimation of  $g_q$ , then running RFTL on vector sequence  $\tilde{\mathbf{g}}_q$  with  $\Phi(x)$  as the regularizer will produce a sequence of point  $\{x_q\}_{q=1}^Q$ ,  $x_q \in \mathcal{K}$ . For  $\{x_q\}_{q=1}^Q$  and any  $\mathbf{y} \in \mathcal{K}$ , we have*

$$\begin{aligned} & \sum_{q=1}^Q \mathbb{E} [\langle \mathbf{g}_q, \mathbf{y} - x_q \rangle] \\ & \leq \eta \sum_{q=1}^Q \mathbb{E} \left[ \|\tilde{\mathbf{g}}_q\|_{\Phi, x_q, *}^2 \right] + \frac{\Phi(\mathbf{y}) - \Phi(x_1)}{\eta} \end{aligned}$$

## 2.2. Ellipsoid Gradient Estimator

Ellipsoid gradient estimator is proposed in (Abernethy et al., 2008), where the authors use it along with the tool from Section 2.1 to design an  $\tilde{O}(\sqrt{T})$  regret algorithm for bandit linear optimization. For a continuous function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  and an invertible matrix  $\mathbf{H} \in \mathbb{R}^{d \times d}$ , we define the  $\mathbf{H}$ -smoothed version of  $f$ .

**Definition 2.4** ( $\mathbf{H}$ -smoothed function). For function  $f(\mathbf{x}): \mathbb{R}^d \rightarrow \mathbb{R}$  and invertible matrix  $\mathbf{H} \in \mathbb{R}^{d \times d}$ , we call  $f^{\mathbf{H}}(\mathbf{x})$  an  $\mathbf{H}$ -smoothed version of  $f(\mathbf{x})$ , where

$$f^{\mathbf{H}}(\mathbf{x}) = \mathbb{E}_{\mathbf{v} \sim \mathbb{B}_d} [f(\mathbf{x} + \mathbf{H}\mathbf{v})].$$

Here  $\mathbf{v} \sim \mathbb{B}_d$  means that  $\mathbf{v}$  is sampled from the unit ball  $\mathbb{B}_d$  uniformly at random.

There is a surprising fact that there is an unbiased estimator of  $\nabla f^{\mathbf{H}}(\mathbf{x})$  for any  $\mathbf{x}$ , and the estimator uses only one query to the value oracle of  $f$ .

**Lemma 2.5** (Ellipsoid estimator (Abernethy et al., 2008)). *Let  $\mathbf{H} \in \mathbb{R}^{d \times d}$  be an invertible matrix,  $f(\mathbf{x}): \mathbb{R}^d \rightarrow \mathbb{R}$  be an arbitrary function. Then*

$$\nabla f^{\mathbf{H}}(\mathbf{x}) = d \mathbb{E}_{\mathbf{v} \sim \mathbb{S}_{d-1}} [f(\mathbf{x} + \mathbf{H}\mathbf{v}) \mathbf{H}^{-1} \mathbf{v}].$$

Here  $\mathbf{v} \sim \mathbb{S}_{d-1}$  means that  $\mathbf{v}$  is sampled from the  $(d-1)$ -dimensional unit sphere  $\mathbb{S}_{d-1}$  uniformly at random.

For linear  $f$ ,  $f^{\mathbf{H}}(\mathbf{x}) = f(\mathbf{x})$ , so Lemma 2.5 gives a one-sample unbiased estimator of the gradient of the linear function. The ellipsoid gradient estimator is usually used along with RFTL with a self-concordant regularizer  $\Phi$  of  $\mathcal{K}$ . When the invertible matrix  $\mathbf{H}$  is set to be  $(\nabla^2 \Phi(\mathbf{x}))^{-1/2}$  and  $\mathbf{x} \in \text{int}(\mathcal{K})$ , the sampled action  $\mathbf{x} + \mathbf{H}\mathbf{v}$  is located in the surface of a so-called **Dikin ellipsoid** centered at  $\mathbf{x}$ , i.e.  $\{\mathbf{x}' \mid \|\mathbf{x}' - \mathbf{x}\|_{\Phi, \mathbf{x}} \leq 1\}$ . The fact that Dikin ellipsoid is entirely contained in  $\mathcal{K}$  is useful for reducing regret.

## 2.3. Non-oblivious Techniques for Monotone DR-Submodular Maximization

The non-oblivious technique was first proposed to improve the approximation ratio of the solution returned by a local search algorithm. The idea is to run a local search on an auxiliary function rather than the original objective, and the local optima of the auxiliary function have a higher approximation ratio, thus the search algorithm will return a better solution.

In monotone submodular maximization literature, Filmus & Ward (2014) improved the approximation ratio of the greedy algorithm to  $1 - 1/e$  using the non-oblivious technique. Zhang et al. (2022a) generalized this result to the continuous DR-submodular maximization problem, improving the approximation ratio of projected gradient ascent to  $1 - 1/e$ . For a monotone DR-submodular function  $f(\mathbf{x})$  satisfying  $f(\mathbf{0}) = 0$ , they consider following auxiliary function,

$$F(\mathbf{x}) = \int_0^1 \frac{e^{z-1}}{z} f(z \cdot \mathbf{x}) dz. \quad (1)$$

We need the following lemma about the auxiliary function proved in their paper.

**Lemma 2.6** (Auxiliary function (Zhang et al., 2022a)). *Let  $f$  be a monotone DR-submodular function defined on  $\mathcal{X}$  and  $f(\mathbf{0}) = 0$ ,  $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ . Let  $F$  be defined as (1). Then*

$$\nabla F(\mathbf{x}) = \int_0^1 e^{z-1} \nabla f(z \cdot \mathbf{x}) dz \quad (2)$$

and the following inequality holds,

$$\langle \mathbf{y} - \mathbf{x}, \nabla F(\mathbf{x}) \rangle \geq (1 - 1/e) f(\mathbf{y}) - f(\mathbf{x}).$$

**Algorithm 1** BanditMLSM( $\eta, L, \Phi$ )

**Input:** block size  $L$ , block number  $Q = T/L$ , learning rate  $\eta$ , self-concordant barrier  $\Phi$

```

1: initiate  $\mathbf{x}_1 \in \text{int}(\mathcal{K})$  such that  $\nabla\Phi(\mathbf{x}_1) = \mathbf{0}$ 
2: for  $q = 1, 2, \dots, Q$  do
3:   Draw  $t_q \sim \text{Unif}\{(q-1)L+1, (q-1)L+2, \dots, qL\}$ 
4:   for  $t = (q-1)L+1, (q-1)L+2, \dots, qL$  do
5:     if  $t = t_q$  then
6:        $\mathbf{H}_q = (\nabla^2\Phi(\mathbf{x}_q))^{-1/2}$ 
7:       sample  $z_q$  from  $\mathbf{Z}$  where  $P(\mathbf{Z} \leq z) = \int_0^z \frac{e^{u-1}}{1-e^{-1}} \mathbb{I}[u \in [0, 1]] du$ 
8:       draw  $\mathbf{v}_q \sim \mathbb{S}_{d-1}$ 
9:       draw  $\mathbf{u}_q$  from  $\{\mathbf{0}, \mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d\}$  with probability:  $\Pr(\mathbf{u}_q = \mathbf{0}) = \frac{1}{2}$ ,  $\Pr(\mathbf{u}_q = \mathbf{e}_i) = \frac{1}{2d}$ 
10:      play  $\mathbf{y}_{t_q} = z_q \cdot \mathbf{x}_q + z_q \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle \mathbf{u}_q$ 
11:      Set  $\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q)$  as (5)
12:       $\tilde{\nabla} \bar{F}_q(\mathbf{x}_q) \leftarrow d \cdot \tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) \mathbf{H}_q^{-1} \mathbf{v}_q$ 
13:       $\mathbf{x}_{q+1} \leftarrow \underset{\mathbf{x} \in \mathcal{K}}{\text{argmin}} \sum_{s=1}^q \langle -\eta \tilde{\nabla} F_s(\mathbf{x}_s), \mathbf{x} \rangle + \Phi(\mathbf{x})$ 
14:    else
15:      play  $\mathbf{y}_t = \mathbf{x}_q$ 
16:    end if
17:  end for
18: end for

```

### 3. Bandit Monotone Multi-linear DR-Submodular Maximization

In this section, we present our algorithm BanditMLSM for BMMDSM. The pseudo-code is shown in Algorithm 1. For some technical reason we will explain later, we divide the whole  $T$  rounds into  $Q$  equal-size blocks, and each block has  $L$  consecutive rounds. Here  $Q$  and  $L$  are to be determined later,  $L = T/Q$ . without loss of generality, we assume both  $L$  and  $Q$  are integers. We define the average function  $\bar{f}_q(\mathbf{x})$  of each block,

$$\bar{f}_q(\mathbf{x}) = \frac{1}{L} \sum_{t=(q-1)L+1}^{qL} f_t(\mathbf{x}). \quad (3)$$

Let  $\bar{F}_q(\mathbf{x})$  be the auxiliary function of  $\bar{f}_q(\mathbf{x})$ ,

$$\bar{F}_q(\mathbf{x}) = \int_0^1 \frac{e^{z-1}}{zL} \sum_{t=(q-1)L+1}^{qL} f_t(z \cdot \mathbf{x}) dz. \quad (4)$$

In high level, BanditMLSM runs RFTL with a self-concordant regularizer  $\Phi(\mathbf{x})$  on the vector sequence  $\{\nabla \bar{F}_q(\mathbf{x}_q)\}_{q=1}^Q$  and controls the regret w.r.t. the linear function sequence  $\{l_q\}_{q=1}^Q$  where  $l_q(\mathbf{u}) := \langle \mathbf{u}, \nabla \bar{F}_q(\mathbf{x}_q) \rangle$ . Now the question is how to estimate  $\nabla \bar{F}_q(\mathbf{x}_q)$ . Recall

Lemma 2.5, we can estimate  $\nabla \bar{F}_q(\mathbf{x}_q) = \nabla l_q(\mathbf{0})$  with the ellipsoid estimator by querying one function value of  $l_q(\mathbf{u})$ . That is, we fix an invertible matrix  $\mathbf{H}_q = (\nabla^2 \Phi(\mathbf{x}_q))^{-1/2}$ , sample a random direction  $\mathbf{v}_q$  in the  $(d-1)$ -dimensional sphere, then query  $l_q(\mathbf{H}_q \mathbf{v}_q)$ , and return  $\tilde{\nabla} \bar{F}_q(\mathbf{x}_q) := d \cdot l_q(\mathbf{H}_q \mathbf{v}_q) \mathbf{H}_q^{-1} \mathbf{v}_q$  as the estimate.

The problem here is that we cannot query  $l_q$  directly. The algorithm can only query the function value of  $f_t$  by playing the corresponding action in round  $t$ . We construct the unbiased estimator of  $l_q(\mathbf{H}_q \mathbf{v}_q)$  as follows. First, we sample a uniformly random  $t_q \in [(q-1)L+1, qL] \cap \mathbb{Z}$  and sample  $z_q$  from the distribution  $Z$  where  $\Pr(Z \leq z) = \int_0^z \frac{e^{u-1}}{1-e^{-1}} \mathbb{I}[u \in [0, 1]] du$ . Then we pick a vector  $\mathbf{u}_q$  from the set  $\{\mathbf{0}, \mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d\}$  following the distribution:  $\Pr(\mathbf{u}_q = \mathbf{0}) = \frac{1}{2}$ ,  $\Pr(\mathbf{u}_q = \mathbf{e}_i) = \frac{1}{2d}$ . Then we play  $\mathbf{y}_{t_q} := z_q \mathbf{x}_q + z_q \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle \mathbf{u}_q$  in round  $t_q$  to obtain the feedback  $f_{t_q}(\mathbf{y}_{t_q})$ . We replace  $l_q(\mathbf{H}_q \mathbf{v}_q)$  with an estimate

$$\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) := \begin{cases} -2(1-1/e) \frac{d}{z_q} \cdot f_{t_q}(\mathbf{y}_{t_q}) & \text{if } \mathbf{u}_q = \mathbf{0}, \\ 2(1-1/e) \frac{d}{z_q} \cdot f_{t_q}(\mathbf{y}_{t_q}) & \text{if } \mathbf{u}_q \neq \mathbf{0}. \end{cases} \quad (5)$$

If  $z_q = 0$ , we define  $\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) := 0$ . The following Lemma shows that  $\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q)$  is an unbiased estimator of  $l_q(\mathbf{H}_q \mathbf{v}_q)$ . Its proof is deferred to Appendix B. In the rounds other than  $t_q$  in block  $q$ , we play  $\mathbf{y}_t := \mathbf{x}_q$  output by RFTL at the end of  $(q-1)$ -th block to exploit the regret bound of RFTL.

**Lemma 3.1.** *Let  $\mathcal{H}_{q-1}$  be the history of the algorithm in the first  $q$  blocks, that is, the realization of  $t_s, z_s, \mathbf{v}_s, \mathbf{u}_s, \forall s \leq q$ . Then  $\mathbb{E}[l_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q] = l_q(\mathbf{H}_q \mathbf{v}_q)$ .*

So  $\tilde{\nabla} \bar{F}_q(\mathbf{x})$  is actually defined as

$$\tilde{\nabla} \bar{F}_q(\mathbf{x}) := d \cdot \tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) \mathbf{H}_q^{-1} \mathbf{v}_q. \quad (6)$$

We show that  $\tilde{\nabla} \bar{F}_q(\mathbf{x}_q)$  is an unbiased estimator of  $\nabla \bar{F}_q(\mathbf{x}_q)$ , and its dual local norm is  $O(d^4)$  in the following lemma. The proof is deferred to Appendix B.

**Lemma 3.2.** *The following properties hold for  $\tilde{\nabla} \bar{F}_q(\mathbf{x}_q)$ :*

- (i)  $\mathbb{E}[\tilde{\nabla} \bar{F}_q(\mathbf{x}_q) \mid \mathcal{H}_{q-1}] = \nabla \bar{F}_q(\mathbf{x}_q)$ ,
- (ii)  $\mathbb{E}[\|\tilde{\nabla} \bar{F}_q(\mathbf{x}_q)\|_{\mathbf{x}_q, *}^2 \mid \mathcal{H}_{q-1}] \leq 4(1-1/e)^2 L_1^2 D^2 d^4$ .

Note that, to estimate  $\nabla \bar{F}_q(\mathbf{x}_q)$ , we must sample an action that is far from  $\mathbf{x}_q$ . This means we cannot do exploration and exploitation at the same time, which is different from the linear bandit. This is the reason why previous works on bandit submodular maximization and our work divide

rounds into blocks. We need to do the exploitation in most of the rounds of a block to maintain the regret bound.

Now recall [Theorem 2.3](#), running RFTL with a self-concordant barrier of  $\mathcal{K}$  will generate a series of action  $\{x_q\}_{q=1}^Q$ , which has low regret w.r.t. the linear function sequence  $\langle \cdot, \nabla \bar{F}_q(x_q) \rangle$ .  $\bar{F}_q$  is the auxiliary function of the block average of  $\{f_t\}_{t=1}^T$ . In block  $q$ , our algorithm plays  $\mathbf{y}_t = x_q$  most of the time. Intuitively, the regret of  $\mathbf{y}_t$  w.r.t. function sequence  $\langle \cdot, \nabla F_t(\mathbf{y}_t) \rangle$  is low, where  $F_t$  is the auxiliary function of  $f_t$ . By [Lemma 2.6](#), we can bound the  $(1 - 1/e)$ -regret of `BanditMLSM`. The proof of [Theorem 3.3](#) is deferred to [Appendix B](#).

**Theorem 3.3.** *Set  $\eta = d^{-4}T^{-2/3}$ ,  $L = d^{-2}T^{1/3}$  in [Algorithm 1](#), if  $\Phi$  is a  $\nu$ -self-concordant barrier of  $\mathcal{K}$ , then the expected  $(1 - 1/e)$ -regret of [Algorithm 1](#) can be bounded as*

$$\mathcal{R}_{1-1/e}(T) \leq O(\nu d^{4/3} T^{2/3} \log T).$$

**About the computational complexity** The computational cost mainly comes from two tasks: (1) Calculating the inverse and square root of the Hessian matrix of the regularizer; (2) Minimizing the convex function over a convex body. These tasks are commonly performed, so `BanditMLSM` can be implemented efficiently.

## 4. Bandit DR-submodular Maximization

Combining RFTL with a self-concordant barrier and non-oblivious technique, we can also improve the result of the general bandit DR-submodular maximization problem where the online reward functions are not required to be multi-linear functions. Due to the space limitation, the algorithmic details and the proof are deferred to the [Appendix C](#). Here we only give the regret bound of our algorithm.

**Theorem 4.1.** *If there is a  $\nu$ -self-concordant barrier of  $\mathcal{K}$ . Then there is an algorithm that attains the following regret upper bound in any  $(\mathcal{K}, \mathcal{F}_{DS})$ -bandit instance:*

$$\mathcal{R}_{1-1/e}(T) \leq O(\nu d^{1/2} T^{3/4} \log T).$$

## 5. A Continuous Approach for Submodular Full-Bandit

In this section, we show reductions from two selected submodular full-bandit problems to the bandit multi-linear DR-submodular maximization problem. All proofs in this section are deferred to [Appendix E](#) due to space limitations.

### 5.1. Reduction Framework

A natural reduction for our task is to consider the multi-linear extension of the submodular function. That is, we consider the multi-linear extension of each submodular set

---

### Algorithm 2 `MLSMWrapper`( $\eta, L, \Phi, \text{EXT}$ )

---

**Input:** learning rate  $\eta$ , block size  $L$ , self-concordant barrier  $\Phi$ , an extension mapping `EXT`

- 1: **for**  $t = 1, 2, \dots, T$  **do**
  - 2:   Get  $\mathbf{y}_t$  from `BanditMLSM4PS`( $\eta, L, \Phi$ )
  - 3:   Sample  $S_t$  from distribution `EXT`( $\mathbf{y}_t$ )
  - 4:   Play  $S_t$  and feed  $g_t(S_t)$  back to `BanditMLSM4PS`( $\eta, L, \Phi$ )
  - 5: **end for**
- 

function, running the `BanditMLSM` on the function sequence of the multi-linear extensions. If we could estimate the function value of the multi-linear extension unbiasedly by using only one query to the corresponding discrete submodular function, then we would complete the reduction successfully. This idea is already considered in ([Zhang et al., 2019](#)). However, it does not work in the full-bandit setting here. The main reason is that the definition of multi-linear extension uses information of the values of the submodular set function on all subsets, including those not satisfying the constraint. This makes it impossible to find an unbiased estimator for the multi-linear extension under bandit feedback setting. To address this problem, [Zhang et al. \(2019\)](#) consider a relaxed responsive bandit model, where they allow the algorithm to query the function value of an infeasible action and gain zero reward. Through this relaxation, they prove a  $O(T^{8/9})$   $(1 - 1/e)$ -regret upper bound for bandit submodular maximization with a matroid constraint. We do not make this relaxation and consider the original full-bandit model, that is, the algorithm must play a feasible action each round.

Assume we want to transform a  $(\mathcal{S}, \mathcal{G})$ -bandit to a bandit multi-linear DR-submodular maximization instance, where  $\mathcal{S}$  is a finite set and we use  $g_t \in \mathcal{G}$  to denote the online reward function. The central component of our reduction framework is a mapping from a product of standard simplexes, denoted as  $\mathcal{K}$ , to the set of all distributions over  $\mathcal{S}$ , denoted as  $\Delta(\mathcal{S})$ . The  $d$ -dimensional standard simplex is a set  $\{(x_1, x_2, \dots, x_d) \mid x_1 + \dots + x_d \leq 1, x_i \geq 0, \forall i\}$ .

We denote the extension mapping as `EXT` :  $\mathcal{K} \rightarrow \Delta(\mathcal{S})$ . The dimension  $d$  of the set  $\mathcal{K}$  varies with different  $\mathcal{S}$  and  $\mathcal{G}$ . The extension mapping naturally defines an extension of any function  $g \in \mathcal{G}$ , that is,  $f(\mathbf{x}) = \mathbb{E}_{S \in \text{EXT}(\mathbf{x})}[g(S)]$ . This extension has a good property, if we sample an element  $S \in \mathcal{S}$  according to the distribution `EXT`( $\mathbf{x}$ ), then  $g(S)$  is an unbiased estimator of  $f(\mathbf{x})$ . The idea is to run `BanditMLSM` on such extensions  $\{f_t\}_{t=1}^T$  of the online functions sequence  $\{g_t\}_{t=1}^T$ . When we received an action  $\mathbf{y}_t$  from `BanditMLSM`, we sample an action  $S_t \in \mathcal{S}$  from `EXT`( $\mathbf{y}_t$ ), and feed  $g_t(S_t)$  back to `BanditMLSM`. However, if we replace the  $f_t(\mathbf{y}_{t_q})$  with  $g_t(S_t)$  in (5), the

estimator  $\tilde{l}(\mathbf{H}_q \mathbf{v}_q)$  can be unbounded when  $z_q$  is very small which makes the regret uncontrollable. Fortunately, when  $\mathcal{K}$  is a product of simplexes, we can slightly modify BanditMLSM to address this problem. We denote the modified algorithm as BanditMLSM4PS. In this algorithm, we use another estimator to substitute (5) when  $z_q < \frac{1}{2}$ . That is, we draw  $\mathbf{u}_q \in \{e_1, \dots, e_d\}$  uniformly at random. Then we let  $\mathbf{y}_{t_q} = z_q \mathbf{x}_q$  or  $\mathbf{y}_{t_q} = z_q \mathbf{x}_q + \frac{1}{2} \mathbf{u}_q$  with equal probability. The estimator is set to be  $\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) := \begin{cases} -4(1-1/e)d\langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle f_{t_q}(\mathbf{y}_{t_q}) & \text{if } \mathbf{y}_{t_q} = z_q \mathbf{x}_q, \\ 4(1-1/e)d\langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle f_{t_q}(\mathbf{y}_{t_q}) & \text{if } \mathbf{y}_{t_q} = z_q \mathbf{x}_q + \frac{1}{2} \mathbf{u}_q. \end{cases}$

To show the estimator is feasible, we need to prove that  $z_q \mathbf{x}_q + \frac{1}{2} \mathbf{u}_q \in \mathcal{K}$  such that the value  $f_{t_q}(z_q \mathbf{x}_q + \frac{1}{2} \mathbf{u}_q)$  can be observed in bandit feedback model. Consider the simplex to which the basis vector  $\mathbf{u}_q$  belongs, without loss of generality, we assume that  $\mathbf{u}_q = e_1$  and  $e_1, e_2, \dots, e_{d_1}$  form the basis of the simplex. Then  $x_1 + \dots + x_{d_1} \leq 1$ , which means  $z_q \sum_{i=1}^{d_1} x_i < \frac{1}{2}$ , therefore  $\frac{1}{2} + z_q \sum_{i=1}^{d_1} x_i \leq 1$ ,  $z_q \mathbf{x}_q + \frac{1}{2} \mathbf{u}_q \in \mathcal{K}$ .

The reduction algorithm is shown in Algorithm 2 and the detailed pseudo-code of BanditMLSM4PS can be found in Algorithm 4 of Appendix E. For product simplexes, we give an  $O(d)$ -self-concordant barrier in Appendix D. To obtain the regret guarantee, we need to make sure that the extension induced by the extension mapping satisfies the assumption BanditMLSM4PS requires. Formally, we prove the following lemma.

**Lemma 5.1.** *For a finite set  $\mathcal{S}$ , and a function family  $\mathcal{G} \subseteq \mathcal{S}^{\mathbb{R}^+}$ , where  $\mathcal{S}^{\mathbb{R}^+}$  is the set of all functions that map element in  $\mathcal{S}$  to  $\mathbb{R}^+$ . If there is an extension mapping  $EXT : \mathcal{K} \rightarrow \Delta(\mathcal{S})$  satisfying following conditions:*

1.  $\mathcal{K} \subseteq \mathbb{R}^d$  is a product of standard simplexes.
2. For any  $g \in \mathcal{G}$ ,  $f(\mathbf{x}) = \mathbb{E}_{S \in EXT(\mathbf{x})}[g(S)]$  is a multi-linear, monotone, DR-submodular function, and  $f$  is  $L_1$ -lipschitz continuous,  $f(\mathbf{0}) = 0$ .
3. For any  $S \in \mathcal{S}$ , there exist  $\mathbf{x} \in \mathcal{K}$  such that  $EXT(\mathbf{x}) = \mathbf{1}_S$ . Where  $\mathbf{1}_S$  assign probability 1 to  $S$  and 0 to other elements of  $\mathcal{S}$ .

then the algorithm MLSMWrapper attains expected  $(1-1/e)$ -regret

$$\mathcal{R}_{1-1/e}(T) \leq O\left(d^{5/3} T^{2/3} \log(T)\right)$$

on  $(\mathcal{S}, \mathcal{G})$ -bandit.

## 5.2. Bandit Monotone Submodular Maximization with Partition Matroid Constraint

We consider a  $(\mathcal{S}_{PM}, \mathcal{G}_{MS})$ -bandit this section, here  $\mathcal{S}_{PM}$  is a partition matroid, and  $\mathcal{G}_{MS}$  is the family of monotone

submodular set function. We assume the functions in  $\mathcal{G}_{MS}$  take value 0 on the empty set.

**Definition 5.2** (Partition Matroid). Let  $G$  be a finite ground set. A set system  $\mathcal{S} \subseteq 2^G$  is called a partition matroid if there exist  $K > 0$  and positive integers  $r_1, r_2, \dots, r_K$  such that  $G$  can be partitioned into  $K$  subsets  $G = \bigcup_{k=1}^K G_k$ , and  $\mathcal{S} = \{A \mid A \in 2^G \text{ and } |A \cap G_k| \leq r_k \forall k\}$ .

By Lemma 5.1, all we need is to find an appropriate extension mapping. Let  $\Delta_d$  be a  $d$ -dimensional standard simplex, Let  $\mathcal{K} = \prod_{k=1}^K \left( \prod_{i=1}^{r_k} \Delta_{|G_k|}^{k,i} \right)$  be the product of standard simplexes. Here  $\Delta_{|G_k|}^{k,i}$  is a  $|G_k|$ -dimensional standard simplex and  $(k, i)$  is the index of this simplex. Next, we construct an extension mapping  $EXT_{PM} : \mathcal{K} \rightarrow \mathcal{S}_{PM}$ .

For  $\mathbf{x} \in \mathcal{K}$ , write  $\mathbf{x} = (x_{k,i,s})_{(k,i,s) \in \Lambda}$ ,  $\Lambda = \{(k, i, s) \mid 1 \leq k \leq K, 1 \leq i \leq r_k, s \in G_k, k, i \in \mathbb{N}\}$  is the index set.  $x_{k,i,s}$  means the coordinate of the simplex  $\Delta_{|G_k|}^{k,i}$ , and  $\mathbf{x} \in \mathbb{R}_+^{\sum_{k=1}^K r_k |G_k|}$  satisfies  $\sum_{s \in G_k} x_{k,i,s} \leq 1, \forall k, i$ . We see the point in the standard simplex  $\Delta_{|G_k|}^{k,i}$  as a probability distribution over  $G_k \cup \{\circ\}$  where  $\circ \notin G_k$  is an extra element which means no element in  $G_k$  is chosen. We sample elements according to the coordinate of each simplex independently, then  $\mathbf{x}$  can be seen as a probability distribution over the set  $\Omega := \prod_{k=1}^K (G_k \cup \{\circ\})^{r_k}$ , we use pre- $EXT_{PM}(\mathbf{x})$  to denote this distribution on  $\Omega$ . We now define a mapping  $\rho : \Omega \rightarrow \mathcal{S}$  as follows. For  $\omega \in \Omega$ , assume  $\omega$  can be represented as  $\omega = (\omega_{k,i})_{(k,i) \in \Gamma}$ , where  $\omega_{k,i} \in G_k \cup \{\circ\}$  and  $\Gamma = \{(k, i) \mid 1 \leq k \leq K, 1 \leq i \leq r_k, k, i \in \mathbb{N}\}$  is the index set. Then  $\rho(\omega) = \{\omega_{k,i} \mid (k, i) \in \Gamma\} \setminus \{\circ\}$ .

It's easy to check  $\rho(\omega) \in \mathcal{S}$ . Thus, for  $\mathbf{x}$ , we first sample an  $\omega \sim \text{pre-}EXT_{PM}(\mathbf{x})$ , then map the sample to  $\rho(\omega) \in \mathcal{S}$ . This process defines a distribution over  $\mathcal{S}$ . We let this distribution be  $EXT_{PM}(\mathbf{x})$ .

**Lemma 5.3.** *For  $\mathcal{G}_{MS}$ , the extension mapping  $EXT_{PM} : \mathcal{K} \rightarrow \Delta(\mathcal{S}_{PM})$  satisfies the conditions in Lemma 5.1. Moreover,  $\mathcal{K}$  is in a  $\sum_{k=1}^K r_k$  dimensional real vector space. For any  $g \in \mathcal{G}_{MS}$ , the continuous extension  $f(\mathbf{x}) = \mathbb{E}_{S \in EXT_{PM}(\mathbf{x})}[g(S)]$  is  $M \sqrt{\sum_{k=1}^K r_k |G_k|}$ -lipschitz.*

**Corollary 5.4.** *There is an algorithm attaining the expected  $(1-1/e)$ -regret of  $\mathcal{R}_{1-1/e}(T) \leq O\left(\left(\sum_{k=1}^K r_k |G_k|\right)^{5/3} T^{2/3} \log(T)\right)$  on any  $(\mathcal{S}_{PM}, \mathcal{G}_{MS})$ -bandit.*

## 5.3. Bandit Sequential Submodular Maximization

Bandit sequential submodular maximization is first studied in (Niazadeh et al., 2021). It is motivated by online retailing platforms where the platform needs to show its products in sequence. There are many types of customers who have different patience and preference. Rarely customers will



see all the products in the list. They will stop browsing the product after seeing some products according to their patience, and the click probability after a customer sees a set of products is submodular. This situation can be formalized into an  $(\mathcal{S}_{OL}, \mathcal{G}_{SS})$ -bandit. Let  $G$  be the ground set of all products, and the constraint  $\mathcal{S}_{OL}$  is the set of all ordered lists of length  $|G|$  consisting of elements in  $G$ .  $\mathcal{G}_{SS}$  consists of function  $g : \mathcal{S}_{OL} \rightarrow [0, M]$  in this form,

$$g(S) = \sum_{i=1}^{|G|} \lambda_i g_i(\{S_j \mid j \leq i\}),$$

where  $\lambda_i$ 's with  $\lambda_i \geq 0$  are positive weights,  $g_i$ 's are monotone submodular set functions, and  $S_j$  is the  $i$ -th element in the ordered list  $S$ . If we interpret  $g(S)$  as a click probability, then  $M = 1$ .

For technical reasons, we assume that there is a dummy element  $\circ$  in  $G$ , which has 0 marginal gain for all  $g_i$ . That is,  $\forall i, \forall S \subseteq G$ , we have  $g_i(S \cup \{\circ\}) = g_i(S)$ . We denote  $G' = G \setminus \{\circ\}$ . This assumption can be satisfied by adding a non-clickable item that is not related to the products to  $G'$ .

Next, we construct an extension mapping  $\text{EXT}_{SS}$ . Let  $\mathcal{K}$  be the cartesian product of standard simplexes,  $\mathcal{K} = \prod_{i=1}^{|G'|} \Delta_{|G'|}^i$ . We see  $\mathbf{x} \in \mathcal{K}$  as  $|G'|$  probability distributions over  $G$ , the component of  $\mathbf{x} \in \mathcal{K}$  in  $\Delta_{|G'|}^i$  represents the distribution of the  $i$ -th element of the ordered list, all these distributions are independent. Any  $\mathbf{x} \in \mathcal{K}$  can be seen as a distribution over  $\mathcal{S}_{OL}$ . Let this distribution be  $\text{EXT}_{SS}(\mathbf{x})$ .

**Lemma 5.5.** *For  $\mathcal{G}_{SS}$ , the extension mapping  $\text{EXT}_{SS} : \mathcal{K} \rightarrow \Delta(\mathcal{S}_{OL})$  satisfies the conditions in Lemma 5.1. Moreover,  $\mathcal{K}$  is in a  $|G|^2 - |G|$  dimensional real vector space. For any  $g \in \mathcal{G}_{SS}$ , the continuous extension  $f(\mathbf{x}) = \mathbb{E}_{S \in \text{EXT}_{SS}(\mathbf{x})}[g(S)]$  is  $M|G|$ -lipschitz.*

**Corollary 5.6.** *There is an algorithm for attaining the expected  $(1 - 1/e)$ -regret of  $\mathcal{R}_{1-1/e}(T) \leq O((|G|)^{10/3} T^{2/3} \log T)$  on any  $(\mathcal{S}_{OL}, \mathcal{G}_{SS})$ -bandit.*

## 6. Conclusion

In this paper, we propose two bandit algorithms, BanditMLSM for monotone multilinear DR-submodular functions and BanditDRSM for general monotone DR-submodular functions. We then show an approach to design the  $\tilde{O}(T^{2/3}) (1 - 1/e)$ -regret algorithm for two special combinatorial full-bandits submodular maximization problems, that is, reducing the combinatorial bandits to a multilinear DR-submodular bandit.

There are some remaining open problems that need to be studied. Firstly, we notice that  $\tilde{O}(T^{2/3})$ -type regret bounds show up frequently in the submodular bandit literature. However, as far as we know no one has proved or disproved the optimality of this bound, which may be an interesting

and challenging problem. Secondly, we still know less about the relationship between the combinatorial constraint and sublinear regret. For submodular set functions, we show that one can achieve sublinear regret with partition matroid constraint in this paper. However, we conjecture that it does not hold for all matroid constraints. How to characterize such a relationship is also a fascinating open question.

## Acknowledgements

We thank the anonymous reviewers for their suggestions in the presentation of the article. This work was supported in part by the National Natural Science Foundation of China Grants No. 61832003, 62272441.

## References

- Abernethy, J., Hazan, E. E., and Rakhlin, A. Competing in the dark: An efficient algorithm for bandit linear optimization. In *21st Annual Conference on Learning Theory, COLT 2008*, pp. 263–273, 2008.
- Bian, A., Levy, K. Y., Krause, A., and Buhmann, J. M. Continuous dr-submodular maximization: structure and algorithms. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 486–496, 2017a.
- Bian, A. A., Mirzasoleiman, B., Buhmann, J., and Krause, A. Guaranteed non-convex optimization: Submodular maximization over continuous domains. In *Artificial Intelligence and Statistics*, pp. 111–120. PMLR, 2017b.
- Chen, L., Krause, A., and Karbasi, A. Interactive submodular bandit. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 140–151, 2017.
- Chen, L., Hassani, H., and Karbasi, A. Online continuous submodular maximization. In *International Conference on Artificial Intelligence and Statistics*, pp. 1896–1905. PMLR, 2018.
- Chen, W., Hu, W., Li, F., Li, J., Liu, Y., and Lu, P. Combinatorial multi-armed bandit with general reward functions. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pp. 1659–1667, 2016.
- Filmus, Y. and Ward, J. Monotone submodular maximization over a matroid via non-oblivious local search. *SIAM Journal on Computing*, 43(2):514–542, 2014.
- Flaxman, A. D., Kalai, A. T., and McMahan, H. B. Online convex optimization in the bandit setting: gradient descent without a gradient. In *Proceedings of the sixteenth*

- annual ACM-SIAM symposium on Discrete algorithms, pp. 385–394, 2005.
- Foster, D. P. and Rakhlin, A. On submodular contextual bandits. *arXiv preprint arXiv:2112.02165*, 2021.
- Hazan, E. et al. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2(3-4):157–325, 2016.
- Li, S., Kong, F., Tang, K., Li, Q., and Chen, W. Online influence maximization under linear threshold model. *Advances in Neural Information Processing Systems*, 33: 1192–1204, 2020.
- Nesterov, Y. and Nemirovskii, A. *Interior-point polynomial algorithms in convex programming*. SIAM, 1994.
- Niazadeh, R., Roughgarden, T., and Wang, J. R. Optimal algorithms for continuous non-monotone submodular and dr-submodular maximization. *The Journal of Machine Learning Research*, 21(1):4937–4967, 2020.
- Niazadeh, R., Golrezaei, N., Wang, J. R., Susan, F., and Badanidiyuru, A. Online learning via offline greedy algorithms: Applications in market design and optimization. In *Proceedings of the 22nd ACM Conference on Economics and Computation*, pp. 737–738, 2021.
- Nie, G., Agarwal, M., Umrawal, A. K., Aggarwal, V., and Quinn, C. J. An explore-then-commit algorithm for submodular maximization under full-bandit feedback. In *The 38th Conference on Uncertainty in Artificial Intelligence*, 2022.
- Nie, G., Nadew, Y. Y., Zhu, Y., Aggarwal, V., and Quinn, C. J. A framework for adapting offline algorithms to solve combinatorial multi-armed bandit problems with bandit feedback. *arXiv preprint arXiv:2301.13326*, 2023.
- Raut, P. S., Sadeghi, O., and Fazel, M. Online dr-submodular maximization with stochastic cumulative constraints. *arXiv preprint arXiv:2005.14708*, 2020.
- Saha, A. and Tewari, A. Improved regret guarantees for online smooth convex optimization with bandit feedback. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 636–642. JMLR Workshop and Conference Proceedings, 2011.
- Streeter, M. and Golovin, D. An online algorithm for maximizing submodular functions. In *Proceedings of the 21st International Conference on Neural Information Processing Systems*, pp. 1577–1584, 2008.
- Streeter, M., Golovin, D., and Krause, A. Online learning of assignments. In *Proceedings of the 22nd International Conference on Neural Information Processing Systems*, pp. 1794–1802, 2009.
- Thang, N. K. and Srivastav, A. Online non-monotone dr-submodular maximization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 9868–9876, 2021.
- Vaswani, S., Lakshmanan, L., Schmidt, M., et al. Influence maximization with bandits. *arXiv preprint arXiv:1503.00024*, 2015.
- Wang, Q. and Chen, W. Improving regret bounds for combinatorial semi-bandits with probabilistically triggered arms and its applications. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 1161–1171, 2017.
- Wu, Q., Li, Z., Wang, H., Chen, W., and Wang, H. Factorization bandits for online influence maximization. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 636–646, 2019.
- Yue, Y. and Guestrin, C. Linear submodular bandits and their application to diversified retrieval. In *Proceedings of the 24th International Conference on Neural Information Processing Systems*, pp. 2483–2491, 2011.
- Zhang, M., Chen, L., Hassani, H., and Karbasi, A. Online continuous submodular maximization: from full-information to bandit feedback. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pp. 9210–9221, 2019.
- Zhang, Q., Deng, Z., Chen, Z., Hu, H., and Yang, Y. Stochastic continuous submodular maximization: Boosting via non-oblivious function. In *International Conference on Machine Learning*, pp. 26116–26134. PMLR, 2022a.
- Zhang, Q., Deng, Z., Chen, Z., Zhou, K., Hu, H., and Yang, Y. Online learning for non-monotone submodular maximization: From full information to bandit feedback. *arXiv preprint arXiv:2208.07632*, 2022b.
- Zhang, Z., Chen, W., Sun, X., and Zhang, J. Online influence maximization with node-level feedback using standard offline oracles. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 9153–9161, 2022c.

## A. Technical Lemmas

This section provides some technical lemmas that will be used in the proofs of this appendix later.

**Definition A.1** (Minkowski function and Minkowski set). Let  $\mathcal{K}$  be a compact convex set, the Minkowski function  $\pi_{\mathbf{x}} : \mathcal{K} \rightarrow \mathbb{R}$  parameterized by a pole  $\mathbf{x} \in \text{int}(\mathcal{K})$  is defined as  $\pi_{\mathbf{x}}(\mathbf{y}) \triangleq \inf\{t \geq 0 \mid \mathbf{x} + t^{-1}(\mathbf{y} - \mathbf{x}) \in \mathcal{K}\}$ . Given  $\delta \in \mathbb{R}^+$  and  $\mathbf{x}_1 \in \text{int}(\mathcal{K})$ , we define the Minkowski set  $\mathcal{K}_{\gamma, \mathbf{x}_1} \triangleq \{\mathbf{x} \in \mathcal{K} \mid \pi_{\mathbf{x}_1}(\mathbf{x}) \leq (1 + \gamma)^{-1}\}$ .

The following lemma provides an upper bound of the difference between the function value of a self-concordant barrier at two different points.

**Lemma A.2** ((Nesterov & Nemirovskii, 1994)). *Let  $\Phi$  be a  $\nu$ -self-concordant barrier over a compact convex set  $\mathcal{K}$ , then for all  $x, y \in \text{int}(\mathcal{K})$ :*

$$\Phi(y) - \Phi(x) \leq \nu \log \frac{1}{1 - \pi_x(y)}.$$

The following lemma is already proved in (Abernethy et al., 2008), we include the proof for completeness.

**Lemma A.3** ((Abernethy et al., 2008)). *Let  $\mathcal{K}$  be a compact convex set,  $\mathbf{x} \in \text{int}(\mathcal{K})$  with diameter  $D$ ,  $\mathbf{x}^* \in \mathcal{K}$  and  $\hat{\mathbf{x}}^* \triangleq \arg\min_{\mathbf{z} \in \mathcal{K}_{\gamma, \mathbf{x}}} \|\mathbf{z} - \mathbf{x}^*\|$  be the projection of  $\mathbf{x}^*$  onto the Minkowski set  $\mathcal{K}_{\gamma, \mathbf{x}}$ , then*

$$\|\mathbf{x}^* - \hat{\mathbf{x}}^*\| \leq \gamma D$$

*Proof.* Consider the point  $\mathbf{y}$  in the segment  $[\mathbf{x}, \mathbf{x}^*]$  satisfying  $\frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{x}^* - \mathbf{x}\|} = \frac{1}{1 + \gamma}$ . Since  $\mathbf{x} + (1 + \gamma)(\mathbf{y} - \mathbf{x}) = \mathbf{x}^* \in \mathcal{K}$ , we can deduce that  $\mathbf{y} \in \mathcal{K}_{\gamma, \mathbf{x}}$ . Thus,

$$\|\hat{\mathbf{x}}^* - \mathbf{x}^*\| \leq \|\mathbf{y} - \mathbf{x}^*\| = \left(1 - \frac{1}{1 + \gamma}\right) \|\mathbf{x}^* - \mathbf{x}\| \leq \gamma D.$$

□

The following two lemmas show that the average auxiliary functions and the  $\mathbf{H}$ -smoothed functions both inherent good properties of the original online functions. And they will be used later.

**Lemma A.4.** *If  $\forall t \in [(q - 1)L + 1, qL]$ ,  $f_t$  is twice differentiable,  $L_1$ -lipschitz and  $L_2$ -smooth, monotone, DR-submodular, then following holds for the average functions  $\bar{f}_q, \bar{F}_q$ .*

- (i)  $\bar{f}_q$  is  $L_1$ -lipschitz and  $L_2$ -smooth.
- (ii)  $\bar{f}_q$  is a monotone DR-submodular function.
- (iii)  $\bar{F}_q$  is  $\frac{L_2}{e}$ -smooth.
- (iv)  $\bar{F}_q$  is a monotone DR-submodular function.

*Proof.* (i)

$$\begin{aligned} \|\bar{f}_q(\mathbf{x}) - \bar{f}_q(\mathbf{y})\| &= \frac{1}{L} \left\| \sum_{t=(q-1)L+1}^{qL} f_t(\mathbf{x}) - \sum_{t=(q-1)L+1}^{qL} f_t(\mathbf{y}) \right\| \\ &\leq \frac{1}{L} \sum_{t=(q-1)L+1}^{qL} \|f_t(\mathbf{x}) - f_t(\mathbf{y})\| \\ &\leq \frac{1}{L} \sum_{t=(q-1)L+1}^{qL} L_1 \|\mathbf{x} - \mathbf{y}\| = L_1 \|\mathbf{x} - \mathbf{y}\| \end{aligned}$$

$$\begin{aligned}
 \|\nabla \bar{f}_q(\mathbf{x}) - \nabla \bar{f}_q(\mathbf{y})\| &= \frac{1}{L} \left\| \sum_{t=(q-1)L+1}^{qL} \nabla f_t(\mathbf{x}) - \sum_{t=(q-1)L+1}^{qL} \nabla f_t(\mathbf{y}) \right\| \\
 &\leq \frac{1}{L} \sum_{t=(q-1)L+1}^{qL} \|\nabla f_t(\mathbf{x}) - \nabla f_t(\mathbf{y})\| \\
 &\leq \frac{1}{L} \sum_{t=(q-1)L+1}^{qL} L_2 \|\mathbf{x} - \mathbf{y}\| \leq L_2 \|\mathbf{x} - \mathbf{y}\|
 \end{aligned}$$

(ii) For any  $i \in [d]$ ,

$$\frac{\partial \bar{f}_q}{\partial x_i}(\mathbf{x}) = \frac{1}{L} \sum_{t=(q-1)L+1}^{qL} \frac{\partial f_t}{\partial x_i}(\mathbf{x}) \geq 0$$

For any  $i \in [d], j \in [d]$ ,

$$\frac{\partial^2}{\partial x_i \partial x_j} \bar{f}_q(\mathbf{x}) = \frac{1}{L} \sum_{t=(q-1)L+1}^{qL} \frac{\partial^2 f_t}{\partial x_i \partial x_j}(\mathbf{x}) \leq 0$$

Thus  $\bar{f}_q$  is monotone DR-submodular.

(iii)

$$\begin{aligned}
 \|\nabla \bar{F}_q(\mathbf{x}) - \nabla \bar{F}_q(\mathbf{y})\| &= \left\| \nabla \int_0^1 \frac{e^{z-1}}{zL} \sum_{t=(q-1)L+1}^{qL} f_t(z \cdot \mathbf{x}) dz - \nabla \int_0^1 \frac{e^{z-1}}{zL} \sum_{t=(q-1)L+1}^{qL} f_t(z \cdot \mathbf{y}) dz \right\| \\
 &= \left\| \int_0^1 \frac{e^{z-1}}{L} \sum_{t=(q-1)L+1}^{qL} \nabla f_t(z \cdot \mathbf{x}) dz - \int_0^1 \frac{e^{z-1}}{L} \sum_{t=(q-1)L+1}^{qL} \nabla f_t(z \cdot \mathbf{y}) dz \right\| \\
 &\leq \int_0^1 \frac{e^{z-1}}{L} \sum_{t=(q-1)L+1}^{qL} \|\nabla f_t(z \cdot \mathbf{x}) - \nabla f_t(z \cdot \mathbf{y})\| dz \\
 &\leq \int_0^1 \frac{e^{z-1}}{L} \sum_{t=(q-1)L+1}^{qL} L_2 z \|\mathbf{x} - \mathbf{y}\| dz \\
 &= L_2 \int_0^1 z e^{z-1} dz \|\mathbf{x} - \mathbf{y}\| = \frac{L_2}{e} \|\mathbf{x} - \mathbf{y}\|
 \end{aligned}$$

(iv) For any  $i \in [d]$ ,

$$\begin{aligned}
 \frac{\partial}{\partial x_i} \bar{F}_q(\mathbf{x}) &= \frac{\partial}{\partial x_i} \int_0^1 \frac{e^{z-1}}{zL} \sum_{t=(q-1)L+1}^{qL} f_t(z \cdot \mathbf{x}) dz \\
 &= \int_0^1 \frac{e^{z-1}}{L} \sum_{t=(q-1)L+1}^{qL} \frac{\partial}{\partial x_i} f_t(z \cdot \mathbf{x}) dz \\
 &\geq 0
 \end{aligned}$$

For any  $i \in [d], j \in [d]$ ,

$$\frac{\partial}{\partial x_i \partial x_j} \bar{F}_q(\mathbf{x}) = \int_0^1 \frac{z e^{z-1}}{L} \sum_{t=(q-1)L+1}^{qL} \frac{\partial}{\partial x_i \partial x_j} f_t(z \cdot \mathbf{x}) dz$$

$$\leq 0$$

Thus  $\overline{F}_q$  is monotone and DR-submodular. □

**Lemma A.5.** *Following properties hold for  $\mathbf{H}$ -smoothed version of a twice differentiable function  $f(\mathbf{x})$ .*

- (i) *If  $f(\mathbf{x})$  is a monotone DR-submodular function, then for any invertible matrix  $\mathbf{H}$ , its  $\mathbf{H}$ -smoothed version  $f^{\mathbf{H}}(\mathbf{x})$  is a monotone DR-submodular function.*
- (ii) *If  $f(\mathbf{x})$  is  $L_1$ -lipschitz continuous and  $L_2$ -smooth, then  $f^{\mathbf{H}}(\mathbf{x})$  is  $L_1$ -lipschitz continuous and  $L_2$ -smooth.*

*Proof.* (i) By Leibnez integral rule, for any  $i \in [d]$ ,

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}_i} f^{\mathbf{H}}(\mathbf{x}) &= \int_{\mathbf{v} \in \mathbb{B}_d} \frac{1}{\text{Vol}(\mathbb{B}_d)} \frac{\partial}{\partial \mathbf{x}_i} f(\mathbf{x} + \mathbf{H}\mathbf{v}) d\mathbf{v} \\ &\geq 0 \end{aligned}$$

The last inequality is because  $\frac{\partial}{\partial \mathbf{x}_i} f(\mathbf{x} + \mathbf{H}\mathbf{v}) \geq 0$  for any  $i \in [d]$ .

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}_i \mathbf{x}_j} f^{\mathbf{H}}(\mathbf{x}) &= \int_{\mathbf{v} \in \mathbb{B}_d} \frac{1}{\text{Vol}(\mathbb{B}_d)} \frac{\partial}{\partial \mathbf{x}_i \mathbf{x}_j} f(\mathbf{x} + \mathbf{H}\mathbf{v}) d\mathbf{v} \\ &\leq 0 \end{aligned}$$

The last inequality is because  $\frac{\partial}{\partial \mathbf{x}_i \mathbf{x}_j} f(\mathbf{x} + \mathbf{H}\mathbf{v}) \leq 0$  for any  $i, j \in [d]$ .

(ii)

$$\begin{aligned} f^{\mathbf{H}}(\mathbf{x}) - f^{\mathbf{H}}(\mathbf{y}) &= \int_{\mathbf{v} \in \mathbb{B}_d} \frac{1}{\text{Vol}(\mathbb{B}_d)} (f(\mathbf{x} + \mathbf{H}\mathbf{v}) - f(\mathbf{y} + \mathbf{H}\mathbf{v})) d\mathbf{v} \\ &\leq \int_{\mathbf{v} \in \mathbb{B}_d} \frac{1}{\text{Vol}(\mathbb{B}_d)} L_1 \|\mathbf{x} + \mathbf{H}\mathbf{v} - \mathbf{y} - \mathbf{H}\mathbf{v}\| d\mathbf{v} \\ &= L_1 \|\mathbf{x} - \mathbf{y}\| \end{aligned}$$

Thus,  $f^{\mathbf{H}}(\mathbf{x})$  is  $L_1$ -lipschitz continuous.

$$\begin{aligned} \nabla f^{\mathbf{H}}(\mathbf{x}) - \nabla f^{\mathbf{H}}(\mathbf{y}) &= \nabla \int_{\mathbf{v} \in \mathbb{B}_d} \frac{1}{\text{Vol}(\mathbb{B}_d)} (f(\mathbf{x} + \mathbf{H}\mathbf{v}) - f(\mathbf{y} + \mathbf{H}\mathbf{v})) d\mathbf{v} \\ &= \int_{\mathbf{v} \in \mathbb{B}_d} \frac{1}{\text{Vol}(\mathbb{B}_d)} \nabla (f(\mathbf{x} + \mathbf{H}\mathbf{v}) - f(\mathbf{y} + \mathbf{H}\mathbf{v})) d\mathbf{v} \\ &\leq \int_{\mathbf{v} \in \mathbb{B}_d} \frac{1}{\text{Vol}(\mathbb{B}_d)} L_2 \|\mathbf{x} - \mathbf{y}\| d\mathbf{v} \\ &= L_2 \|\mathbf{x} - \mathbf{y}\| \end{aligned}$$

□

## B. Missing Proofs in Section 3

We first show a property of multi-linear functions, which is the key observation of our estimator for the gradient of multi-linear functions.

**Lemma B.1.** *If  $f : \mathcal{K} \rightarrow \mathbb{R}$  is a multi-linear function, where  $\mathcal{K} \subseteq \mathbb{R}^d$ , then for any basis vector  $\mathbf{e}_i, i \in [d]$  and any  $\mathbf{x} \in \mathcal{K}$  and  $\lambda > 0$  satisfying  $\mathbf{x} + \lambda \mathbf{e}_i \in \mathcal{K}$ . We have,*

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}_i} = \frac{f(\mathbf{x} + \lambda \mathbf{e}_i) - f(\mathbf{x})}{\lambda} \tag{7}$$

*Proof.* When we fix the components of  $\mathbf{x}$  except  $x_i$ ,  $f$  is a linear function of  $x_i$ . Thus (7) directly comes from the linearity.  $\square$

**Lemma 3.1.** *Let  $\mathcal{H}_q$  be the history of the algorithm in the first  $q$  blocks, that is, the realization of  $t_s, z_s, \mathbf{v}_s, \mathbf{u}_s, \forall s \leq q$ . Then  $\mathbb{E}[\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q] = l_q(\mathbf{H}_q \mathbf{v}_q)$ .*

*Proof of Lemma 3.1.* When  $z_q > 0$ , we have

$$\begin{aligned} \mathbb{E}[\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q, t_q, z_q] &= \frac{1}{2}(-2(1-1/e)\frac{d}{z_q} \cdot f_{t_q}(z_q \mathbf{x}_q)) + \sum_{i=1}^d \frac{1}{2d}(2(1-1/e)\frac{d}{z_q} \cdot f_{t_q}(z_q \mathbf{x}_q + z_q \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle \mathbf{e}_i)) \\ &= (1-1/e) \sum_{i=1}^d \frac{1}{z_q} (f_{t_q}(z_q \mathbf{x}_q + z_q \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle \mathbf{e}_i) - f_{t_q}(z_q \mathbf{x}_q)) \\ &= (1-1/e) \sum_{i=1}^d \frac{1}{z_q} z_q \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle \frac{\partial f_{t_q}}{\partial x_i}(z_q \cdot \mathbf{x}_q) \\ &= (1-1/e) \sum_{i=1}^d \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle \langle \mathbf{e}_i, \nabla f_{t_q}(z_q \cdot \mathbf{x}_q) \rangle \\ &= (1-1/e) \langle \mathbf{H}_q \mathbf{v}_q, \nabla f_{t_q}(z_q \cdot \mathbf{x}_q) \rangle \end{aligned}$$

Then we take the expectations over  $t_q$  and  $z_q$ , note the value of  $\tilde{l}(\mathbf{H}_q \mathbf{v}_q)$  when  $z_q = 0$  does not affect the result of the integral since  $z_q = 0$  is a zero measured event.

$$\begin{aligned} \mathbb{E}[\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q] &= \sum_{t_q=(q-1)L+1}^{qL} \int_0^1 \Pr(t_q, z_q) \mathbb{E}[\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q, t_q, z_q] dz_q \\ &= \sum_{t_q=(q-1)L+1}^{qL} \int_0^1 \frac{1}{L} \frac{e^{z_q-1}}{1-1/e} (1-1/e) \langle \mathbf{H}_q \mathbf{v}_q, \nabla f_{t_q}(z_q \cdot \mathbf{x}_q) \rangle dz_q \\ &= \left\langle \mathbf{H}_q \mathbf{v}_q, \int_0^1 \frac{e^{z_q-1}}{L} \sum_{t_q=(q-1)L+1}^{qL} \nabla f_{t_q}(z_q \cdot \mathbf{x}_q) dz_q \right\rangle \end{aligned}$$

Since

$$\begin{aligned} \nabla \bar{F}_q(\mathbf{x}_q) &= \nabla \int_0^1 \frac{e^{z-1}}{zL} \sum_{t=(q-1)L+1}^{qL} f_t(z \cdot \mathbf{x}_q) dz \\ &= \int_0^1 \frac{e^{z-1}}{L} \sum_{t=(q-1)L+1}^{qL} \nabla f_t(z \cdot \mathbf{x}_q) dz \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbb{E}[\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q] &= \langle \mathbf{H}_q \mathbf{v}_q, \nabla \bar{F}_q(\mathbf{x}_q) \rangle \\ &= l_q(\mathbf{H}_q \mathbf{v}_q) \end{aligned}$$

$\square$

**Lemma 3.2.** *The following properties hold for  $\tilde{\nabla} \bar{F}_q(\mathbf{x}_q)$*

$$(i) \mathbb{E} \left[ \tilde{\nabla} \bar{F}_q(\mathbf{x}_q) \mid \mathcal{H}_{q-1} \right] = \nabla \bar{F}_q(\mathbf{x}_q)$$

$$(ii) \mathbb{E} \left[ \|\tilde{\nabla} \bar{F}_q(\mathbf{x}_q)\|_{\mathbf{x}_q, *}^2 \mid \mathcal{H}_{q-1} \right] \leq 4(1-1/e)^2 L_1^2 D^2 d^4$$

*Proof of Lemma 3.2.*

(i)

$$\begin{aligned} \mathbb{E} \left[ \widetilde{\nabla} \overline{F}_q(\mathbf{x}_q) \mid \mathcal{H}_{q-1} \right] &= \int_{\mathbf{v}_q \in \mathbb{S}_{d-1}} d \cdot \mathbb{E}[\widetilde{l}_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q] \mathbf{H}_q^{-1} \mathbf{v}_q d\mathbf{v}_q \\ &= \int_{\mathbf{v}_q \in \mathbb{S}_{d-1}} \frac{1}{\text{Vol}(\mathbb{S}_{d-1})} d \cdot l_q(\mathbf{H}_q \mathbf{v}_q) \mathbf{H}_q^{-1} \mathbf{v}_q d\mathbf{v}_q \end{aligned} \quad (8)$$

$$\begin{aligned} &= \mathbb{E}_{\mathbf{v}_q \sim \mathbb{S}_{d-1}} [d \cdot l_q(\mathbf{H}_q \mathbf{v}_q) \mathbf{H}_q^{-1} \mathbf{v}_q] \\ &= \nabla l_q^{\mathbf{H}_q}(\mathbf{0}) \end{aligned} \quad (9)$$

$$\begin{aligned} &= \nabla l_q(\mathbf{0}) \\ &= \nabla \overline{F}_q(\mathbf{x}_q) \end{aligned} \quad (10)$$

(8) is due to Lemma 3.1, (9) is due to Lemma 2.5, (10) is because that  $l_q$  is a linear function.

(ii)

$$\begin{aligned} &\mathbb{E} \left[ \|\widetilde{\nabla} \overline{F}_q(\mathbf{x}_q)\|_{\mathbf{x}_q, * }^2 \mid \mathcal{H}_{q-1} \right] \\ &= \mathbb{E} \left[ \frac{1}{2} (1 - 1/e)^2 (2d^2)^2 \frac{1}{z_q^2} f_{t_q}(z_q \cdot \mathbf{x}_q)^2 \cdot \mathbf{v}_q^T \mathbf{H}_q \Phi(\mathbf{x}_q)^{-1} \mathbf{H}_q^{-1} \mathbf{v}_q \mid \mathcal{H}_{q-1} \right] \\ &\quad + \sum_{i=1}^d \mathbb{E} \left[ \frac{1}{2d} (1 - 1/e)^2 4d^4 \frac{1}{z_q^2} f_{t_q}(z_q \cdot \mathbf{x}_q + z_q \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle \mathbf{e}_i)^2 \mathbf{v}_q^T \mathbf{H}_q \Phi(\mathbf{x}_q) \mathbf{H}_q^{-1} \mathbf{v}_q \mid \mathcal{H}_{q-1} \right] \\ &\leq 2(1 - 1/e)^2 d^4 \frac{L_1^2 z_q^2 \|\mathbf{x}_q\|^2}{z_q^2} \mathbb{E} [\mathbf{v}_q^T \mathbf{H}_q^{-1} \Phi(\mathbf{x}_q)^{-1} \mathbf{H}_q^{-1} \mathbf{v}_q \mid \mathcal{H}_{q-1}] \\ &\quad + 2(1 - 1/e)^2 d^4 \frac{L_1^2 z_q^2 \|\mathbf{x}_q + \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle \mathbf{e}_i\|^2}{z_q^2} \mathbb{E} [\mathbf{v}_q^T \mathbf{H}_q^{-1} \Phi(\mathbf{x}_q)^{-1} \mathbf{H}_q^{-1} \mathbf{v}_q \mid \mathcal{H}_{q-1}] \end{aligned} \quad (11)$$

$$\begin{aligned} &= 4(1 - 1/e)^2 d^4 L_1^2 D^2 \mathbb{E} \left[ \mathbf{v}_q^T (\Phi(\mathbf{x}_q)^{-1/2})^{-1} \Phi(\mathbf{x}_q)^{-1} (\Phi(\mathbf{x}_q)^{-1/2})^{-1} \mathbf{v}_q \mid \mathcal{H}_{q-1} \right] \\ &\leq 4(1 - 1/e)^2 d^4 L_1^2 D^2 \|\mathbf{v}_q\|_2^2 \\ &= 4(1 - 1/e)^2 d^4 L_1^2 D^2 \end{aligned} \quad (12)$$

Inequality (11) is because  $f_{t_q}$  is  $L_1$ -lipschitz continuous and  $\mathbf{x}_q + \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle \mathbf{e}_i$  is in the Dikin ellipsoid  $\{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}_q\|_{\Phi, \mathbf{x}_q} \leq 1\}$ , which is contained in  $\mathcal{K}$ . (12) is because  $\mathbf{v}_q \in \mathbb{S}_{d-1}$ , thus  $\|\mathbf{v}_q\| = 1$ .

□

**Theorem 3.3.** Set  $\eta = d^{-4} T^{-2/3}$ ,  $L = d^{-2} T^{1/3}$ ,  $Q = T/L = d^2 T^{2/3}$  in Algorithm 1, if  $\Phi$  is a  $\nu$ -self-concordant barrier of  $\mathcal{K}$ , then the expected  $(1 - 1/e)$ -regret of Algorithm 1 can be bounded as

$$\mathcal{R}_{1-1/e}(T) \leq (4(1 - 1/e)^2 L_1^2 D^2 + M) d^{4/3} T^{2/3} + (1 - 1/e) L_1 D + \nu d^{4/3} T^{2/3} \log(T)$$

*Proof of Theorem 3.3.* Set  $\widetilde{\mathbf{g}}_q = \widetilde{\nabla} \overline{F}_q(\mathbf{x}_q)$ ,  $\mathbf{g}_q = \nabla \overline{F}_q(\mathbf{x}_q)$  in Theorem 2.3. We have proved  $\mathbb{E} [\widetilde{\nabla} \overline{F}_q(\mathbf{x}_q) \mid \mathcal{H}_{q-1}] = \nabla \overline{F}_q(\mathbf{x}_q)$ . Let  $\hat{\mathbf{x}}^* \triangleq \text{argmin}_{\mathbf{x} \in \mathcal{K}_{\gamma, \mathbf{x}_1}} \|\mathbf{x}^* - \mathbf{x}\|$  be the projection of  $\mathbf{x}^*$  onto the Minkowski set  $\mathcal{K}_{\gamma, \mathbf{x}_1}$  defined in Definition A.1, here the pole is  $\mathbf{x}_1$ , and  $\gamma$  is a parameter to be determined later,  $\mathbf{x}^* = \text{argmin}_{\mathbf{x} \in \mathcal{K}} \sum_{t=1}^T f_t(\mathbf{x})$ . We have

$$\begin{aligned}
 \sum_{q=1}^Q \mathbb{E} [\langle \nabla \bar{F}_q(\mathbf{x}_q), \hat{\mathbf{x}}^* - \mathbf{x}_q \rangle \mid \mathcal{H}_{q-1}] &\leq \eta \sum_{q=1}^Q \mathbb{E} \left[ \|\tilde{\nabla} \bar{F}_q(\mathbf{x}_q)\|_{\mathbf{x}_{q,*}}^2 \mid \mathcal{H}_{q-1} \right] + \frac{\Phi(\hat{\mathbf{x}}^*) - \Phi(\mathbf{x}_1)}{\eta} \\
 &\leq 4(1-1/e)^2 L_1^2 D^2 \eta d^4 Q + \frac{\Phi(\hat{\mathbf{x}}^*) - \Phi(\mathbf{x}_1)}{\eta} \\
 &\leq 4(1-1/e)^2 L_1^2 D^2 \eta d^4 Q + \frac{\nu \log(\frac{1}{1-(1+\gamma)^{-1}})}{\eta}
 \end{aligned} \tag{13}$$

The last inequality is because  $\pi_{\mathbf{x}_1}(\hat{\mathbf{x}}^*) \leq (1+\gamma)^{-1}$  and [Lemma A.2](#). Since  $\bar{f}_q$  is a monotone DR-submodular function by [Lemma A.4](#) and  $\bar{F}_q$  is its auxiliary function, we lower bound the left hand side of (13) by [Lemma 2.6](#),

$$\begin{aligned}
 \sum_{q=1}^Q \mathbb{E} [\langle \nabla \bar{F}_q(\mathbf{x}_q), \hat{\mathbf{x}}^* - \mathbf{x}_q \rangle \mid \mathcal{H}_{q-1}] &\geq \sum_{q=1}^Q \mathbb{E} [((1-1/e)\bar{f}_q(\hat{\mathbf{x}}^*) - \bar{f}_q(\mathbf{x}_q)) \mid \mathcal{H}_{q-1}] \\
 &= \sum_{q=1}^Q \sum_{t=(q-1)L+1}^{qL} \frac{1}{L} \mathbb{E} [((1-1/e)f_t(\hat{\mathbf{x}}^*) - f_t(\mathbf{x}_q)) \mid \mathcal{H}_{q-1}] \\
 &= \frac{1}{L} \sum_{t=1}^T \mathbb{E} \left[ (1-1/e)f_t(\hat{\mathbf{x}}^*) - f_t(\mathbf{x}_{\lceil \frac{t}{L} \rceil}) \mid \mathcal{H}_{\lceil \frac{t}{L} \rceil - 1} \right] \\
 &= \underbrace{\frac{1}{L} \sum_{t=1}^T \mathbb{E} \left[ (1-1/e)f_t(\hat{\mathbf{x}}^*) - (1-1/e)f_t(\mathbf{x}^*) \mid \mathcal{H}_{\lceil \frac{t}{L} \rceil - 1} \right]}_{(A)} \\
 &\quad + \frac{1}{L} \sum_{t=1}^T \mathbb{E} \left[ (1-1/e)f_t(\mathbf{x}^*) - f_t(\mathbf{y}_t) \mid \mathcal{H}_{\lceil \frac{t}{L} \rceil - 1} \right] \\
 &\quad + \underbrace{\frac{1}{L} \sum_{q=1}^Q \mathbb{E} \left[ f_{t_q}(\mathbf{y}_{t_q}) - f_{t_q}(\mathbf{x}_{\lceil \frac{t_q}{L} \rceil}) \mid \mathcal{H}_{q-1} \right]}_{(B)}
 \end{aligned}$$

Since  $|f_t(\hat{\mathbf{x}}^*) - f_t(\mathbf{x}^*)| \leq L_1 \|\hat{\mathbf{x}}^* - \mathbf{x}^*\| \leq L_1 \gamma D$  and  $|f_{t_q}(\mathbf{y}_{t_q}) - f_{t_q}(\mathbf{x}_{\lceil \frac{t_q}{L} \rceil})| \leq M$ , we have

$$|(A)| \leq (1-1/e) \frac{L_1}{L} \gamma DT \quad |(B)| \leq \frac{MQ}{L} \tag{14}$$

Therefore,

$$\begin{aligned}
 \mathbb{E} \left[ \sum_{t=1}^T (1-1/e)f_t(\mathbf{x}^*) - f_t(\mathbf{y}_t) \right] &\leq L \sum_{q=1}^Q \mathbb{E} [\langle \nabla \bar{F}_q(\mathbf{x}_q), \hat{\mathbf{x}}^* - \mathbf{x}_q \rangle \mid \mathcal{H}_{q-1}] - L \times (A) - L \times (B) \\
 &\leq 4(1-1/e)^2 L_1^2 D^2 \eta d^4 T + \frac{\nu L \log(\frac{1}{1-(1+\gamma)^{-1}})}{\eta} + (1-1/e)L_1 \gamma DT + MQ
 \end{aligned}$$

The last inequality is because of (13) and (14). set  $\eta = d^{-8/3} T^{-1/3}$ ,  $L = d^{-4/3} T^{1/3}$ ,  $Q = T/L = d^{4/3} T^{2/3}$ ,  $\gamma = \frac{1}{T}$ , we have,

$$\begin{aligned}
 \mathbb{E} \left[ \sum_{t=1}^T (1-1/e)f_t(\mathbf{x}^*) - f_t(\mathbf{y}_t) \right] &\leq (4(1-1/e)^2 L_1^2 D^2 + M) d^{4/3} T^{2/3} + (1-1/e)L_1 D + \nu d^{4/3} T^{2/3} \log(T+1) \\
 &= O(\nu d^{4/3} T^{2/3} \log(T))
 \end{aligned}$$

□



**Algorithm 3** BanditDRSM( $\eta, \delta, L, \Phi$ )

**Input:** Smoothing radius  $\delta$ , block size  $L$ , block number  $Q = T/L$ , learning rate  $\eta$ , self-concordant barrier  $\Phi$

```

1: initiate  $\mathbf{x}_1 \in \text{int}(\mathcal{K})$  such that  $\nabla\Phi(\mathbf{x}_1) = 0$ 
2: for  $q = 1, 2, \dots, Q$  do
3:   Draw  $t_q \sim \text{Unif}\{(q-1)L+1, (q-1)L+2, \dots, qL\}$ 
4:   for  $t = (q-1)L+1, (q-1)L+2, \dots, qL$  do
5:     if  $t = t_q$  then
6:        $\mathbf{H}_q = (\nabla^2\Phi(\mathbf{x}_q))^{-1/2}$ 
7:       sample  $z_q$  from  $\mathbf{Z}$  where  $P(\mathbf{Z} \leq z) = \int_0^z \frac{e^{u-1}}{1-e^{-1}} \mathbb{I}[u \in [0, 1]] du$ 
8:       draw  $\mathbf{v}_q \sim \mathbb{S}_{d-1}$ 
9:       play  $\mathbf{y}_t = z_q \cdot \mathbf{x}_q + \delta z_q \cdot \mathbf{H}_q \mathbf{v}_q$ 
10:       $\tilde{\nabla}\bar{F}_q(\mathbf{x}_q) \leftarrow (1-1/e) \frac{d}{\delta z_q} f_{t_q}(\mathbf{y}_t) \mathbf{H}_q^{-1} \mathbf{v}_q$ 
11:       $\mathbf{x}_{q+1} \leftarrow \underset{\mathbf{x} \in \mathcal{K}}{\text{argmin}} \sum_{s=1}^q \langle -\eta \tilde{\nabla}\bar{F}_s(\mathbf{x}_s), \mathbf{x} \rangle + \Phi(\mathbf{x})$ 
12:     else
13:       play  $\mathbf{y}_t = \mathbf{x}_q$ 
14:     end if
15:   end for
16: end for

```

### C. Bandit DR-submodular Maximization

In this section we present our algorithm BanditDRSM for general bandit monotone DR-submodular maximization, the pseudocode is shown in Algorithm 3. BanditDRSM is very similar to BanditMLSM, it also divides  $T$  rounds into  $Q$  equal size blocks. We use again  $\bar{f}_q(\mathbf{x})$  and  $\bar{F}_q(\mathbf{x})$  to denote the average function of  $q$ -th block and the auxiliary function of it, defined as (3) and (4). BanditDRSM runs RFTL with self-concordant regularizer on vector sequence  $\{\nabla\bar{F}_q(\mathbf{x}_q)\}_{q=1}^Q$ . Here the difference compared with BanditMLSM is, we cannot find an unbiased estimator for  $\nabla\bar{F}_q(\mathbf{x}_q)$ . We use the ellipsoid estimator directly to estimate  $\nabla\bar{F}_q^{\delta\mathbf{H}_q}(\mathbf{x}_q)$ , the gradient of the  $\delta\mathbf{H}_q$ -smoothed function, here  $\mathbf{H}_q = (\nabla^2\Phi(\mathbf{x}_q))^{-1/2}$  is the same as BanditMLSM,  $\delta$  is a parameter to be determined. Specifically, in block  $q$ , we select a uniform random exploration round  $t_q \in [(q-1)L+1, qL] \cap \mathbb{Z}$ , a random direction  $\mathbf{v}_q \in \mathbb{S}_{d-1}$ ,  $z_q \sim Z$  where  $\Pr(Z \leq z) = \int_0^z \frac{e^{u-1}}{1-e^{-1}} \mathbb{I}[u \in \{0, 1\}] du$ . In round  $t_q$ , we play  $\mathbf{y}_{t_q} = z_q \cdot \mathbf{x}_q + \delta z_q \cdot \mathbf{H}_q \mathbf{v}_q$  and feedback the gradient estimate as follow to RFTL and define the estimator  $\tilde{\nabla}\bar{F}(\mathbf{x}_q)$ .

$$\tilde{\nabla}\bar{F}(\mathbf{x}_q) := (1-1/e) \frac{d}{z_q \delta} \cdot f_{t_q}(\mathbf{y}_{t_q}) \mathbf{H}_q^{-1} \mathbf{v}_q. \quad (15)$$

$\mathbf{y}_{t_q} = z_q(\mathbf{x}_q + \delta\mathbf{H}_q\mathbf{v}_q)$ , if we let  $\delta \leq 1$ , then  $\mathbf{x}_q + \delta\mathbf{H}_q\mathbf{v}_q$  is in the Dikin ellipsoid  $\{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}_q\|_{\Phi, \mathbf{x}_q} \leq 1\}$ , therefore  $\mathbf{x}_q + \delta\mathbf{H}_q^{-1}\mathbf{v}_q \in \mathcal{K}$ . Since  $\mathbf{0} \in \mathcal{K}$ ,  $z_q \in [0, 1]$  and  $\mathcal{K}$  is convex,  $\mathbf{y}_{t_q} \in \mathcal{K}$ . When  $z_q = 0$ , we define  $\tilde{\nabla}\bar{F}(\mathbf{x}_q) := 0$ , since  $\Pr(z_q = 0) = 0$ , the value of  $\tilde{\nabla}\bar{F}(\mathbf{x}_q)$  when  $z_q = 0$  does not matter.

We prove that  $\tilde{\nabla}\bar{F}(\mathbf{x}_q)$  is an unbiased gradient estimator for the  $\delta\mathbf{H}_q$ -smoothed function  $\bar{F}_q^{\delta\mathbf{H}_q}(\mathbf{x}_q)$ . Moreover, the dual local norm of the estimator can be bounded as  $O(\frac{d^2}{\delta^2})$ . To formalize the above arguments, we have the following lemma.

**Lemma C.1.** *Let  $\tilde{\nabla}\bar{F}_q(\mathbf{x}_q)$  be defined as (15). Assume  $f_t$  for  $t \in [(q-1)L+1, qL]$  is  $L_1$ -lipschitz,  $f_t(\mathbf{0}) = 0$ , then the following hold,*

$$(i) \mathbb{E} \left[ \tilde{\nabla}\bar{F}_q(\mathbf{x}_q) \mid \mathcal{H}_{q-1} \right] = \nabla\bar{F}_q^{\delta\mathbf{H}_q}(\mathbf{x}_q).$$

$$(ii) \|\tilde{\nabla}\bar{F}_q(\mathbf{x}_q)\|_{\mathbf{x}_q, *}^2 \leq \frac{(1-e)^2 d^2 L_1^2 D^2}{\delta^2}$$

*Proof.* (i) Let  $\mathbf{H} = \delta \mathbf{H}_q$  in Lemma 2.5, let  $f_{t_q, z_q}(\mathbf{x}) \triangleq f_{t_q}(z_q \cdot \mathbf{x})$ , we have

$$\mathbb{E} \left[ \tilde{\nabla} \bar{F}_q(\mathbf{x}_q) \mid \mathcal{H}_{q-1}, t_q, z_q \right] = \frac{1-1/e}{z_q} \nabla f_{t_q, z_q}^{\delta \mathbf{H}_q}(\mathbf{x}_q)$$

Thus,

$$\begin{aligned} \mathbb{E} \left[ \tilde{\nabla} \bar{F}_q(\mathbf{x}_q) \mid \mathcal{H}_{q-1} \right] &= \sum_{t_q=(q-1)L+1}^{qL} \int_0^1 \Pr(t_q, z_q \mid \mathcal{H}_{q-1}) \frac{1-1/e}{z_q} \nabla f_{t_q, z_q}^{\delta \mathbf{H}_q}(\mathbf{x}_q) dz_q \\ &= \sum_{t_q=(q-1)L+1}^{qL} \int_0^1 \frac{e^{z_q-1}}{(1-1/e)L} \frac{1-1/e}{z_q} \nabla f_{t_q, z_q}^{\delta \mathbf{H}_q}(\mathbf{x}_q) dz_q \\ &= \sum_{t_q=(q-1)L+1}^{qL} \int_0^1 \frac{e^{z_q-1}}{z_q L} \nabla f_{t_q, z_q}^{\delta \mathbf{H}_q}(\mathbf{x}_q) dz_q \\ &= \nabla \bar{F}_q^{\delta \mathbf{H}_q}(\mathbf{x}_q) \end{aligned}$$

(ii)

$$\begin{aligned} \|\tilde{\nabla} \bar{F}_q(\mathbf{x}_q)\|_{\mathbf{x}_q, *}^2 &= (1-1/e)^2 \frac{d^2}{\delta^2 z_q^2} f_{t_q}^2(z_q \cdot \mathbf{x}_q + \delta z_q \cdot \mathbf{H}_q \mathbf{v}_q) \mathbf{v}_q^T \mathbf{H}_q^{-1} (\nabla^2 \Phi(\mathbf{x}_q))^{-1} \mathbf{H}_q^{-1} \mathbf{v}_q \\ &\leq (1-1/e)^2 \frac{d^2}{\delta^2 z_q^2} z_q^2 L_1^2 \|\mathbf{x}_q + \delta \mathbf{H}_q \mathbf{v}_q\|^2 \|\mathbf{v}_q\|^2 \\ &\leq \frac{(1-e)^2 d^2 L_1^2 D^2}{\delta^2} \end{aligned}$$

The first inequality is because  $f_{t_q}$  is  $L_1$ -lipschitz continuous. □

Intuitively, we can control the regret of  $\{\mathbf{x}_q\}_{q=1}^Q$  w.r.t. the linear function sequence  $\{\langle \cdot, \nabla \bar{F}_q^{\delta \mathbf{H}_q}(\mathbf{x}_q) \rangle\}_{q=1}^Q$  by using Theorem 2.3. In Lemma A.5, We proved that  $\bar{f}_q^{\delta \mathbf{H}_q}$  is also DR-submodular. Since  $\bar{F}_q^{\delta \mathbf{H}_q}(\mathbf{x}_q)$  is the auxiliary function of  $\bar{f}_q^{\delta \mathbf{H}_q}$ , this allows us to control the  $(1-1/e)$ -regret of  $\{\mathbf{x}_q\}$  w.r.t.  $\{\bar{f}_q^{\delta \mathbf{H}_q}\}$  by using Lemma 2.6. A key observation here is  $\|\bar{f}_q^{\delta \mathbf{H}_q} - \bar{f}_q\|_\infty \leq O(\delta^2)$  assuming the online functions are smooth, which means we can bound the  $(1-1/e)$ -regret of  $\{\mathbf{x}_q\}$  w.r.t.  $\{\bar{f}_q\}$  in term of the  $(1-1/e)$ -regret w.r.t.  $\{\bar{f}_q^{\delta \mathbf{H}_q}\}$  with an extra  $O(\delta^2)$  additive term. Previous works (Zhang et al., 2019; Niazadeh et al., 2021) use the FKM estimator proposed in (Flaxman et al., 2005), where the sample sphere is fixed (which can be seen as a special case of the ellipsoid estimator when  $\mathbf{H}_q = I$ ), to prevent the sample action jump out  $\mathcal{K}$ , they must run their algorithm on a smaller interior  $\mathcal{K}_\delta$  which is  $\delta$ -far from  $\partial \mathcal{K}$ . So this only guarantees the regret competing with the point in  $\mathcal{K}_\delta$ , this adds an  $O(\delta)$  term to the overall regret, which is bigger than  $O(\delta^2)$  since the  $\delta$  is set to  $o(1)$  latter.

With this improved gradient estimator and non-oblivious technique, we prove a  $\tilde{O}(T^{3/4})$   $(1-1/e)$ -regret of BanditDRSM.

**Theorem 4.1** (restatement). *Set  $\eta = D^{-2} d^{-1} T^{-1/2}$ ,  $\delta = d^{1/4} T^{-1/8}$ ,  $L = d^{-1/2} T^{1/4}$ ,  $Q = T/L = d^{1/2} T^{3/4}$  in Algorithm 3. If  $\Phi$  is a  $\nu$ -self concordant function of  $\mathcal{K}$ , then the expected  $(1-1/e)$ -regret of Algorithm 3 can be bounded as*

$$\mathcal{R}_{1-1/e}(T) \leq O(\nu d^{1/2} T^{3/4} \log(T))$$

*Proof of Theorem 4.1.* Let  $\hat{\mathbf{x}}^* = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}_{\gamma, \mathbf{x}_1}} \|\mathbf{x}^* - \mathbf{x}\|$ , where  $\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x} \in \mathcal{K}} \sum_{t=1}^T f_t(\mathbf{x}^*)$ . Let  $\mathbf{g}_q = \nabla \bar{F}_q^{\delta \mathbf{H}_q}$ ,  $\tilde{\mathbf{g}}_q = \tilde{\nabla} \bar{F}_q(\mathbf{x}_q)$ ,  $\mathbf{y} = \hat{\mathbf{x}}^*$  in Theorem 2.3. Since we proved  $\tilde{\nabla} \bar{F}_q(\mathbf{x}_q)$  is an unbiased estimate of  $\nabla \bar{F}_q^{\delta \mathbf{H}_q}(\mathbf{x}_q)$  in Lemma C.1, we have

$$\sum_{q=1}^Q \mathbb{E} \left[ \langle \nabla \bar{F}_q^{\delta \mathbf{H}_q}(\mathbf{x}_q), \hat{\mathbf{x}}^* - \mathbf{x}_q \rangle \mid \mathcal{H}_{q-1} \right] \leq \eta \sum_{q=1}^Q \mathbb{E} \left[ \|\tilde{\nabla} \bar{F}_q(\mathbf{x}_q)\|_{\Phi, \mathbf{x}_q, *}^2 \mid \mathcal{H}_{q-1} \right] + \frac{\Phi(\hat{\mathbf{x}}^*) - \Phi(\mathbf{x}_1)}{\eta}$$

$$\begin{aligned}
 &\leq \eta \sum_{q=1}^Q \frac{(1-e)^2 d^2 L_1^2 D^2}{\delta^2} + \frac{\Phi(\hat{\mathbf{x}}^*) - \Phi(\mathbf{x}_1)}{\eta} \\
 &\leq \frac{(1-e)^2 \eta Q d^2 L_1^2 D^2}{\delta^2} + \frac{\nu \log\left(\frac{1}{1-(1+\gamma)^{-1}}\right)}{\eta}
 \end{aligned}$$

Since  $\nabla \bar{F}_q^{\delta \mathbf{H}_q}(\mathbf{x})$  is monotone DR-submodular due to [Lemma A.5](#), and it is the auxiliary function of  $\bar{f}_q^{\delta \mathbf{H}_q}(\mathbf{x})$ , by [Lemma 2.6](#), we have

$$\sum_{q=1}^Q \mathbb{E} \left[ \langle \nabla \bar{F}_q^{\delta \mathbf{H}_q}(\mathbf{x}_q), \hat{\mathbf{x}}^* - \mathbf{x}_q \rangle \mid \mathcal{H}_{q-1} \right] \geq \sum_{q=1}^Q \mathbb{E} \left[ (1-1/e) \bar{f}_q^{\delta \mathbf{H}_q}(\hat{\mathbf{x}}^*) - \bar{f}_q^{\delta \mathbf{H}_q}(\mathbf{x}_q) \mid \mathcal{H}_{q-1} \right]$$

The RHS can be further decomposed into several terms,

$$\begin{aligned}
 &\sum_{q=1}^Q \mathbb{E} \left[ (1-1/e) \bar{f}_q^{\delta \mathbf{H}_q}(\hat{\mathbf{x}}^*) - \bar{f}_q^{\delta \mathbf{H}_q}(\mathbf{x}_q) \mid \mathcal{H}_{q-1} \right] \\
 &= \underbrace{\sum_{q=1}^Q \mathbb{E} \left[ (1-1/e) \bar{f}_q^{\delta \mathbf{H}_q}(\hat{\mathbf{x}}^*) - (1-1/e) \bar{f}_q^{\delta \mathbf{H}_q}(\mathbf{x}^*) \mid \mathcal{H}_{q-1} \right]}_{(A)} + \underbrace{\sum_{q=1}^Q \mathbb{E} \left[ (1-1/e) \bar{f}_q^{\delta \mathbf{H}_q}(\mathbf{x}^*) - (1-1/e) \bar{f}_q(\mathbf{x}^*) \mid \mathcal{H}_{q-1} \right]}_{(B)} \\
 &\quad + \underbrace{\sum_{q=1}^Q \mathbb{E} \left[ (1-1/e) \bar{f}_q(\mathbf{x}^*) - \bar{f}_q(\mathbf{x}_q) \mid \mathcal{H}_{q-1} \right]}_{(C)} + \underbrace{\sum_{q=1}^Q \mathbb{E} \left[ \bar{f}_q(\mathbf{x}_q) - \bar{f}_q^{\delta \mathbf{H}_q}(\mathbf{x}_q) \mid \mathcal{H}_{q-1} \right]}_{(C)}
 \end{aligned} \tag{16}$$

**Bounding (A):** Since  $f_t(\mathbf{x})$  is  $L_1$ -lipschitz continuous for any  $t$ ,  $\bar{f}_q$  is also  $L_1$ -lipschitz continuous by [Lemma A.4](#), thus  $\bar{f}_q^{\delta \mathbf{H}_q}$  is  $L_1$ -lipschitz continuous by [Lemma A.5](#). Since  $\|\hat{\mathbf{x}}^* - \mathbf{x}^*\| \leq \gamma D$  by [Lemma A.3](#),

$$\begin{aligned}
 \sum_{q=1}^Q \mathbb{E} \left[ (1-1/e) \bar{f}_q^{\delta \mathbf{H}_q}(\hat{\mathbf{x}}^*) - (1-1/e) \bar{f}_q^{\delta \mathbf{H}_q}(\mathbf{x}^*) \mid \mathcal{H}_{q-1} \right] &\geq - \sum_{q=1}^Q (1-1/e) \mathbb{E} \left[ |\bar{f}_q^{\delta \mathbf{H}_q}(\hat{\mathbf{x}}^*) - \bar{f}_q^{\delta \mathbf{H}_q}(\mathbf{x}^*)| \mid \mathcal{H}_{q-1} \right] \\
 &\geq - \sum_{q=1}^Q (1-1/e) L_1 \gamma D = -(1-1/e) L_1 \gamma D Q
 \end{aligned} \tag{17}$$

**Bounding (B):** Since  $f_t(\mathbf{x})$  is  $L_2$ -smooth for any  $t$ , by [Lemma A.4](#) and [Lemma A.5](#),  $\bar{f}_q^{\delta \mathbf{H}_q}$  is  $L_2$ -smooth. Thus,

$$\begin{aligned}
 \bar{f}_q^{\delta \mathbf{H}_q}(\mathbf{x}^*) - \bar{f}_q(\mathbf{x}^*) &= \frac{1}{\text{Vol}(\mathbb{B}_d)} \int_{\mathbf{v} \in \mathbb{B}_d} \bar{f}_q(\mathbf{x}^* + \delta \mathbf{H}_q \mathbf{v}) - \bar{f}_q(\mathbf{x}^*) d\mathbf{v} \\
 &\geq \frac{1}{\text{Vol}(\mathbb{B}_d)} \int_{\mathbf{v} \in \mathbb{B}_d} \langle \nabla \bar{f}_q(\mathbf{x}^*), \delta \mathbf{H}_q \mathbf{v} \rangle - \frac{L_2}{2} \|\delta \mathbf{H}_q \mathbf{v}\|^2 d\mathbf{v} \\
 &= \frac{1}{\text{Vol}(\mathbb{B}_d)} \left\langle \nabla \bar{f}_q(\mathbf{x}^*), \delta \mathbf{H}_q \int_{\mathbf{v} \in \mathbb{B}_d} \mathbf{v} d\mathbf{v} \right\rangle - \frac{1}{\text{Vol}(\mathbb{B}_d)} \int_{\mathbf{v} \in \mathbb{B}_d} \frac{L_2}{2} \|\delta \mathbf{H}_q \mathbf{v}\|^2 d\mathbf{v} \\
 &\geq - \frac{1}{\text{Vol}(\mathbb{B}_d)} \int_{\mathbf{v} \in \mathbb{B}_d} \frac{L_2}{2} \delta^2 D^2 d\mathbf{v} \\
 &\geq - \frac{L_2 \delta^2 D^2}{2}
 \end{aligned}$$

Therefore,

$$\sum_{q=1}^Q \mathbb{E} \left[ (1-1/e) \bar{f}_q^{\delta \mathbf{H}_q}(\hat{\mathbf{x}}^*) - (1-1/e) \bar{f}_q(\hat{\mathbf{x}}^*) \mid \mathcal{H}_{q-1} \right] \geq - \frac{(1-1/e) L_2 \delta^2 D^2 Q}{2} \tag{18}$$

**Bounding (C):** Similarly,

$$\begin{aligned}\bar{f}_q(\mathbf{x}_q) - \bar{f}_q^{\delta \mathbf{H}_q}(\mathbf{x}_q) &= \frac{1}{\text{Vol}(\mathbb{B}_d)} \int_{\mathbf{v} \in \mathbb{B}_d} \bar{f}_q(\mathbf{x}_q) - \bar{f}_q(\mathbf{x}_q + \delta \mathbf{H}_q \mathbf{v}) d\mathbf{v} \\ &\geq \frac{1}{\text{Vol}(\mathbb{B}_d)} \int_{\mathbf{v} \in \mathbb{B}_d} \langle \nabla \bar{f}_q(\mathbf{x}_q), \delta \mathbf{H}_q \mathbf{v} \rangle - \frac{L_2}{2} \|\delta \mathbf{H}_q \mathbf{v}\|^2 d\mathbf{v} \\ &\geq -\frac{L_2 \delta^2 D^2}{2}\end{aligned}$$

Therefore,

$$\sum_{q=1}^Q \mathbb{E} \left[ \bar{f}_q(\mathbf{x}_q) - \bar{f}_q^{\delta \mathbf{H}_q}(\mathbf{x}_q) \mid \mathcal{H}_{q-1} \right] \geq -\frac{L_2 \delta^2 D^2 Q}{2} \quad (19)$$

Put (17),(18),(19) in (16) and rearrange it,

$$\begin{aligned}&\sum_{q=1}^Q \mathbb{E} \left[ (1 - 1/e) \bar{f}_q(\mathbf{x}^*) - \bar{f}_q(\mathbf{x}_q) \mid \mathcal{H}_{q-1} \right] \\ &\leq \sum_{q=1}^Q \mathbb{E} \left[ \langle \nabla \bar{F}_q^{\delta \mathbf{H}_q}(\mathbf{x}_q), \hat{\mathbf{x}}^* - \mathbf{x}_q \rangle \mid \mathcal{H}_{q-1} \right] + (1 - 1/e) L_1 \gamma D Q + \frac{(1 - 1/e) L_2 \delta^2 D^2 Q}{2} + \frac{L_2 \delta^2 D^2 Q}{2} \\ &\leq \frac{(1 - e)^2 \eta Q d^2 L_1^2 D^2}{\delta^2} + \frac{\nu \log\left(\frac{1}{1 - (1 + \gamma)^{-1}}\right)}{\eta} + (1 - 1/e) L_1 \gamma D Q + \frac{(2 - 1/e) L_2 \delta^2 D^2 Q}{2}\end{aligned}$$

Then we bound the expected regret,

$$\begin{aligned}\mathcal{R}_{1-1/e}(T) &= \sum_{t=1}^T \mathbb{E} \left[ (1 - 1/e) f_t(\mathbf{x}^*) - f_t(\mathbf{y}_t) \mid \mathcal{H}_{\lceil \frac{t}{L} \rceil - 1} \right] \\ &= \sum_{t=1}^T \mathbb{E} \left[ (1 - 1/e) f_t(\mathbf{x}^*) - f_t(\mathbf{x}_{\lceil \frac{t}{L} \rceil}) \mid \mathcal{H}_{\lceil \frac{t}{L} \rceil - 1} \right] + \sum_{q=1}^Q \mathbb{E} \left[ f_{t_q}(\mathbf{x}_q) - f_{t_q}(\mathbf{y}_{t_q}) \mid \mathcal{H}_{q-1} \right] \\ &\leq L \sum_{q=1}^Q \mathbb{E} \left[ (1 - 1/e) \bar{f}_q(\mathbf{x}^*) - \bar{f}_q(\mathbf{x}_q) \mid \mathcal{H}_{q-1} \right] + M Q \\ &\leq \frac{(1 - e)^2 \eta L Q d^2 L_1^2 D^2}{\delta^2} + \frac{\nu L \log\left(\frac{1}{1 - (1 + \gamma)^{-1}}\right)}{\eta} + (1 - 1/e) L_1 \gamma D L Q + \frac{(2 - 1/e) L_2 \delta^2 D^2 L Q}{2} + M Q \\ &= \frac{(1 - e)^2 \eta d^2 L_1^2 D^2 T}{\delta^2} + \frac{\nu \log\left(\frac{1}{1 - (1 + \gamma)^{-1}}\right) L}{\eta} + (1 - 1/e) L_1 \gamma D T + \frac{(2 - 1/e) L_2 \delta^2 D^2 T}{2} + M Q\end{aligned}$$

Set  $\eta = D^{-2} d^{-1} T^{-1/2}$ ,  $\delta = D^{-1/2} d^{1/4} T^{-1/8}$ ,  $L = D^{-1} d^{-1/2} T^{1/4}$ ,  $Q = T/L = D d^{1/2} T^{3/4}$ ,  $\gamma = \frac{1}{T}$

$$\begin{aligned}\mathcal{R}_{1-1/e}(T) &\leq (1 - e)^2 L_1^2 D d^{1/2} T^{3/4} + \nu D d^{1/2} T^{3/4} \log(T + 1) + (1 - 1/e) L_1 D \\ &\quad + \frac{(2 - 1/e) L_2 \delta^2 D d^{1/2} T^{3/4}}{2} + M D d^{1/2} T^{3/4} \\ &= O(\nu d^{1/2} T^{3/4} \log(T))\end{aligned}$$

□

The idea of using a self-concordant regularizer RFTL on smooth online functions is motivated by (Saha & Tewari, 2011). Where the authors studied the bandit convex optimization problem, and they find that RFTL with the self-concordant regularizer works well when the convex functions are smooth. We find this idea also works here in the bandit DR-submodular maximization problem.

## D. Self-Concordant Barrier of Product Simplexes

In this section, we give a self-concordant barrier for the product simplex, which is a cartesian product of several simplexes. Let  $\mathcal{K}$  be the product of  $n$  simplexes, and their dimensions are  $d_1, d_2, \dots, d_n$  respectively. We write  $\mathcal{K}$  as

$$\mathcal{K} = \prod_{i=1}^n \Delta_{d_i}.$$

For  $\mathbf{x} \in \mathcal{K}$ , we represent it as  $\mathbf{x} = (x_{1,1}, x_{1,2}, \dots, x_{1,d_1}, x_{2,1}, \dots, x_{2,d_2}, \dots, x_{n,d_n})$ .  $\mathbf{x} \in \mathcal{K}$  iff

$$\begin{cases} x_{i,j} \geq 0, & \forall 1 \leq i \leq n \text{ and } 1 \leq j \leq d_i \\ \sum_{j=1}^{d_i} x_{i,j} \leq 1, & \forall 1 \leq i \leq n \end{cases}$$

Define the function  $\Phi : \text{int}(\mathcal{K}) \rightarrow \mathbb{R}$ ,

$$\Phi(\mathbf{x}) = - \sum_{i=1}^n \log(1 - \vec{1}_{d_i}^T \cdot \mathbf{x}_i) - \sum_{i=1}^n \sum_{j=1}^{d_i} \log(x_{i,j}).$$

Here  $\vec{1}_{d_i} = \underbrace{(1, 1, \dots, 1)}_{d_i}^T$ ,  $\mathbf{x}_i = (x_{i,1}, x_{i,2}, \dots, x_{i,d_i})$ . We prove that  $\Phi$  is a  $n$ -self-concordant barrier of  $\mathcal{K}$ .

**Lemma D.1.**  $\Phi(\mathbf{x})$  is a  $\sum_{i=1}^n (d_i + 1)$ -self-concordant barrier of  $\mathcal{K}$ .

*Proof.* It's easy to see that  $\Phi(\mathbf{x})$  is three-times continuously differentiable and approaches infinity along any sequence of points approaching the boundary of  $\mathcal{K}$ . We first calculate the gradient and the hessian matrix of  $\Phi$ .

$$\begin{aligned} \frac{\partial \Phi}{\partial x_{i,j}}(\mathbf{x}) &= - \sum_{i=1}^n \frac{\partial \log(1 - \vec{1}_{d_i}^T \cdot \mathbf{x}_i)}{\partial x_{i,j}} - \sum_{i=1}^n \sum_{j=1}^{d_i} \frac{\partial \log(x_{i,j})}{\partial x_{i,j}} \\ &= \frac{1}{1 - \vec{1}^T \cdot \mathbf{x}_i} - \frac{1}{x_{i,j}} \end{aligned}$$

$$\begin{aligned} \frac{\partial^2 \Phi(\mathbf{x})}{\partial x_{i_1, j_1} \partial x_{i_2, j_2}}(\mathbf{x}) &= \frac{\partial(1 - \vec{1}^T \cdot \mathbf{x}_{i_1})^{-1}}{\partial x_{i_1, j_1}} - \frac{\partial x_{i_2, j_2}^{-1}}{\partial x_{i_1, j_1}} \\ &= \frac{1}{(1 - \vec{1}^T \cdot \mathbf{x}_{i_1})^2} \mathbb{I}[i_1 = i_2] + \frac{1}{x_{i_1, j_1}^2} \mathbb{I}[i_1 = i_2, j_1 = j_2] \end{aligned}$$

For any direction  $\mathbf{h} = (h_{1,1}, \dots, h_{1,d_1}, h_{2,1}, \dots, h_{n,d_n})^T$ ,

$$\begin{aligned} \mathbf{h}^T \nabla^2 \Phi(\mathbf{x}) \mathbf{h} &= \sum_{i_1=1}^n \sum_{j_1=1}^{d_{i_1}} \sum_{i_2=1}^n \sum_{j_2=1}^{d_{i_2}} h_{i_1, j_1} h_{i_2, j_2} \frac{\partial^2 \Phi(\mathbf{x})}{\partial x_{i_1, j_1} \partial x_{i_2, j_2}}(\mathbf{x}) \\ &= \sum_{i_1=1}^n \sum_{j_1=1}^{d_{i_1}} \sum_{i_2=1}^n \sum_{j_2=1}^{d_{i_2}} \left( \frac{h_{i_1, j_1} h_{i_2, j_2}}{(1 - \vec{1}^T \cdot \mathbf{x}_{i_1})^2} \mathbb{I}[i_1 = i_2] + \frac{h_{i_1, j_1} h_{i_2, j_2}}{x_{i_1, j_1}^2} \mathbb{I}[i_1 = i_2, j_1 = j_2] \right) \\ &= \sum_{i=1}^n \frac{(\sum_{j=1}^{d_i} h_{i,j})^2}{(1 - \vec{1}^T \cdot \mathbf{x}_{i_1})^2} + \sum_{i=1}^n \sum_{j=1}^{d_i} \frac{h_{i,j}^2}{x_{i,j}^2} \geq 0 \end{aligned}$$

Therefore,  $\Phi$  is convex. Next, we check the condition 2 of Definition 2.1. Let  $\mathbf{h} = (h_{1,1}, h_{1,d_1}, h_{2,1}, \dots, h_{2,d_2}, \dots, h_{n,d_n})$ ,  $\mathbf{h}_i = (h_{i,1}, \dots, h_{i,d_i})$ . Then,

$$\begin{aligned}
 & \nabla^3 \Phi(\mathbf{x})[\mathbf{h}, \mathbf{h}, \mathbf{h}] \\
 &= \frac{\partial^3}{\partial t_1 \partial t_2 \partial t_3} \left( - \sum_{i=1}^n \log(1 - \bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{x}_i - t_1 \bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{h}_i - t_2 \bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{h}_i - t_3 \bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{h}_i) - \sum_{i=1}^n \sum_{j=1}^{d_i} \log(x_{i,j} + (t_1 + t_2 + t_3)h_{i,j}) \right) \Big|_{t_1=t_2=t_3=0} \\
 &= \sum_{i=1}^n \frac{2(\bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{h}_i)^3}{(1 - \bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{x}_i - t_1 \bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{h}_i - t_2 \bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{h}_i - t_3 \bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{h}_i)^3} - \sum_{i=1}^n \sum_{j=1}^{d_i} \frac{2h_{i,j}^3}{(x_{i,j} + (t_1 + t_2 + t_3)h_{i,j})^3} \Big|_{t_1=t_2=t_3=0} \\
 &= \sum_{i=1}^n \frac{2(\bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{h}_i)^3}{(1 - \bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{x}_i)^3} - \sum_{i=1}^n \sum_{j=1}^{d_i} \frac{2h_{i,j}^3}{x_{i,j}^3}
 \end{aligned}$$

We check the first inequality in the condition 2 of Definition 2.1.

$$\begin{aligned}
 2(\nabla^2 \Phi(\mathbf{x})[\mathbf{h}, \mathbf{h}])^{3/2} &= 2 \left( \sum_{i=1}^n \frac{(\sum_{j=1}^{d_i} h_{i,j})^2}{(1 - \bar{\mathbf{1}}^T \cdot \mathbf{x}_{i_1})^2} + \sum_{i=1}^n \sum_{j=1}^{d_i} \frac{h_{i,j}^2}{x_{i,j}^2} \right)^{3/2} \\
 &= 2 \left( \sum_{i=1}^n \frac{(\sum_{j=1}^{d_i} h_{i,j})^2}{(1 - \bar{\mathbf{1}}^T \cdot \mathbf{x}_{i_1})^2} + \sum_{i=1}^n \sum_{j=1}^{d_i} \frac{h_{i,j}^2}{x_{i,j}^2} \right) \left( \sum_{i=1}^n \frac{(\sum_{j=1}^{d_i} h_{i,j})^2}{(1 - \bar{\mathbf{1}}^T \cdot \mathbf{x}_{i_1})^2} + \sum_{i=1}^n \sum_{j=1}^{d_i} \frac{h_{i,j}^2}{x_{i,j}^2} \right)^{1/2} \\
 &= 2 \left( \sum_{i=1}^n \frac{(\sum_{j=1}^{d_i} h_{i,j})^2}{(1 - \bar{\mathbf{1}}^T \cdot \mathbf{x}_{i_1})^2} \left( \sum_{i=1}^n \frac{(\sum_{j=1}^{d_i} h_{i,j})^2}{(1 - \bar{\mathbf{1}}^T \cdot \mathbf{x}_{i_1})^2} + \sum_{i=1}^n \sum_{j=1}^{d_i} \frac{h_{i,j}^2}{x_{i,j}^2} \right)^{1/2} \right. \\
 &\quad \left. + \sum_{i=1}^n \sum_{j=1}^{d_i} \frac{h_{i,j}^2}{x_{i,j}^2} \left( \sum_{i=1}^n \frac{(\sum_{j=1}^{d_i} h_{i,j})^2}{(1 - \bar{\mathbf{1}}^T \cdot \mathbf{x}_{i_1})^2} + \sum_{i=1}^n \sum_{j=1}^{d_i} \frac{h_{i,j}^2}{x_{i,j}^2} \right)^{1/2} \right) \\
 &\geq 2 \left( \sum_{i=1}^n \left| \frac{(\sum_{j=1}^{d_i} h_{i,j})^3}{(1 - \bar{\mathbf{1}}^T \cdot \mathbf{x}_{i_1})^3} \right| + \sum_{i=1}^n \sum_{j=1}^{d_i} \left| \frac{h_{i,j}^3}{x_{i,j}^3} \right| \right) \\
 &\geq \left| \sum_{i=1}^n \frac{2(\bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{h}_i)^3}{(1 - \bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{x}_i)^3} - \sum_{i=1}^n \sum_{j=1}^{d_i} \frac{2h_{i,j}^3}{x_{i,j}^3} \right| = |\nabla^3 \Phi(\mathbf{x})[\mathbf{h}, \mathbf{h}, \mathbf{h}]|
 \end{aligned}$$

Then we check the inequality between  $\nabla \Phi(\mathbf{x})[\mathbf{h}]$  and  $\nabla^2 \Phi(\mathbf{x})[\mathbf{h}, \mathbf{h}]$ .

$$\begin{aligned}
 |\nabla \Phi(\mathbf{x})[\mathbf{h}]| &= |\mathbf{h}^T \nabla \Phi(\mathbf{x})| \\
 &\leq \sum_{i=1}^n \left| \frac{\bar{\mathbf{1}}_{d_i}^T \cdot \mathbf{h}_i}{1 - \bar{\mathbf{1}}^T \cdot \mathbf{x}_i} \right| + \sum_{i=1}^n \sum_{j=1}^{d_i} \left| \frac{h_{i,j}}{x_{i,j}} \right| \\
 &\leq \sqrt{\sum_{i=1}^n (d_i + 1)} \left( \sum_{i=1}^n \frac{(\sum_{j=1}^{d_i} h_{i,j})^2}{(1 - \bar{\mathbf{1}}^T \cdot \mathbf{x}_{i_1})^2} + \sum_{i=1}^n \sum_{j=1}^{d_i} \frac{h_{i,j}^2}{x_{i,j}^2} \right)^{1/2} \\
 &= \left( \sum_{i=1}^n (d_i + 1) \right)^{1/2} (\nabla^2 \Phi(\mathbf{x})[\mathbf{h}, \mathbf{h}])^{1/2}
 \end{aligned}$$

Therefore,  $\Phi(\mathbf{x})$  is a  $(\sum_{i=1}^n (d_i + 1))$ -self-concordant barrier of  $\mathcal{K}$ .  $\square$

**Algorithm 4** BanditMLSM4PS( $\eta, L, \Phi$ )

**Input:** block size  $L$ , block number  $Q = T/L$ , learning rate  $\eta$ , potential function  $\Phi$ 

```

1: initiate  $\mathbf{x}_1 \in \text{int}(\mathcal{K})$  such that  $\nabla\Phi(\mathbf{x}_1) = \mathbf{0}$ 
2: for  $q = 1, 2, \dots, Q$  do
3:   Draw  $t_q \sim \text{Unif}\{(q-1)L+1, (q-1)L+2, \dots, qL\}$ 
4:   for  $t = (q-1)L+1, (q-1)L+2, \dots, qL$  do
5:     if  $t = t_q$  then
6:        $\mathbf{H}_q = (\nabla^2\Phi(\mathbf{x}_q))^{-1/2}$ 
7:       sample  $z_q$  from  $\mathbf{Z}$  where  $P(\mathbf{Z} < z) = \int_0^z \frac{e^{u-1}}{1-e^{-1}} \mathbb{I}[u \in [0, 1]] du$ 
8:       draw  $\mathbf{v}_q \sim \mathbb{S}_{d-1}$ 
9:       if  $z_q \geq \frac{1}{2}$  then
10:        draw  $\mathbf{u}_q$  from  $\{\mathbf{0}, \mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d\}$  with probability:  $\Pr(\mathbf{u}_q = \mathbf{0}) = \frac{1}{2}, \Pr(\mathbf{u}_q = \mathbf{e}_i) = \frac{1}{2d}$ 
11:         $\mathbf{y}_{t_q} \leftarrow z_q \cdot \mathbf{x}_q + z_q \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle \mathbf{u}_q$ 
12:         $\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) \leftarrow \begin{cases} -2(1-1/e) \frac{d}{z_q} \cdot f_{t_q}(\mathbf{y}_{t_q}) & \text{if } \mathbf{u}_q = \mathbf{0}, \\ 2(1-1/e) \frac{d}{z_q} \cdot f_{t_q}(\mathbf{y}_{t_q}) & \text{if } \mathbf{u}_q \neq \mathbf{0}. \end{cases}$ 
13:       else
14:        draw  $\mathbf{u}_q$  from  $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d\}$  uniformly at random
15:        let  $\mathbf{y}_{t_q} = z_q \mathbf{x}_q + \frac{1}{2} \mathbf{u}_q$  or  $\mathbf{y}_{t_q} = z_q \mathbf{x}_q$  with equal probability.
16:        play  $\mathbf{y}_{t_q}$  and observe the feedback  $f_{t_q}(\mathbf{y}_{t_q})$ 
17:         $\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) \leftarrow \begin{cases} -4(1-1/e)d \cdot \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle f_{t_q}(\mathbf{y}_{t_q}) & \text{if } \mathbf{y}_{t_q} = z_q \mathbf{x}_q, \\ 4(1-1/e)d \cdot \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle f_{t_q}(\mathbf{y}_{t_q}) & \text{if } \mathbf{y}_{t_q} = z_q \mathbf{x}_q + \frac{1}{2} \mathbf{u}_q. \end{cases}$ 
18:       end if
19:        $\tilde{\nabla} F_q(\mathbf{x}_q) \leftarrow d \cdot \tilde{l}_q(\mathbf{H}_q \mathbf{v}_q) \mathbf{H}_q^{-1} \mathbf{v}_q$ 
20:        $\mathbf{x}_{q+1} \leftarrow \underset{\mathbf{x} \in \mathcal{K}}{\text{argmin}} \sum_{s=1}^q \langle -\eta \tilde{\nabla} F_s(\mathbf{x}_s), \mathbf{x} \rangle + \Phi(\mathbf{x})$ 
21:     else
22:        $\mathbf{y}_t \leftarrow \mathbf{x}_q$ ,
23:       sample  $S_t$  from  $\text{EXT}(\mathbf{y}_t)$  and play  $S_t$ .
24:     end if
25:   end for
26: end for

```

## E. Missing Proofs in Section 5

The detailed pseudo-code of BanditMLSM4PS is shown in Algorithm 4, the only difference between BanditMLSM4PS and BanditMLSM is the line 9 to line 17 in Algorithm 4.

### E.1. Proof of Lemma 5.1

In Algorithm 2, the algorithm MLSMWrapper feeds  $g_{t_q}(S_{t_q})$  back to BanditMLSM4PS to replace the value  $f_{t_q}(\mathbf{y}_{t_q})$ . Therefore MLSMWrapper are actually using a new estimator for  $l_q(\mathbf{H}_q \mathbf{v}_q)$ , we denote the new estimator  $\tilde{l}'(\mathbf{H}_q \mathbf{v}_q)$ . If  $z_q \geq \frac{1}{2}$ :

$$\tilde{l}'(\mathbf{H}_q \mathbf{v}_q) := \begin{cases} -2(1-1/e) \frac{d}{z_q} \cdot g_{t_q}(S_{t_q}) & \text{if } \mathbf{u}_q = \mathbf{0}, \\ 2(1-1/e) \frac{d}{z_q} \cdot g_{t_q}(S_{t_q}) & \text{if } \mathbf{u}_q \neq \mathbf{0}. \end{cases}$$

else:

$$\tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q) := \begin{cases} -4(1-1/e)d \cdot \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle g_{t_q}(S_{t_q}) & \text{if } \mathbf{y}_{t_q} = z_q \mathbf{x}_q, \\ 4(1-1/e)d \cdot \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle g_{t_q}(S_{t_q}) & \text{if } \mathbf{y}_{t_q} = z_q \mathbf{x}_q + \frac{1}{2} \mathbf{u}_q. \end{cases} \quad (20)$$

We first show  $\tilde{l}(\mathbf{H}_q \mathbf{v}_q)$  is an unbiased estimator of  $l(\mathbf{H}_q \mathbf{v}_q)$ .

**Lemma E.1.** *The estimator  $\tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q)$  is an unbiased estimator for  $l_q(\mathbf{H}_q \mathbf{v}_q)$ , that is,*

$$\mathbb{E} \left[ \tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_q, \mathbf{v}_q \right] = l_q(\mathbf{H}_q \mathbf{v}_q)$$

*Proof.* Condition on  $\mathcal{H}_{q-1}, \mathbf{v}_q, t_q, z_q, \mathbf{u}_q$ . If  $z_q \geq \frac{1}{2}$  and  $\mathbf{u}_q = 0$ ,

$$\begin{aligned} \mathbb{E} \left[ \tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q, t_q, z_q, \mathbf{u}_q \right] &= -2(1-1/e) \frac{d}{z_q} \mathbb{E}[g_{t_q}(S_{t_q}) \mid \mathcal{H}_{q-1}, \mathbf{v}_q, t_q, z_q, \mathbf{u}_q] \\ &= -2(1-1/e) \frac{d}{z_q} f_{t_q}(z_q \mathbf{x}_q). \end{aligned}$$

The last equality is because that  $S_{t_q} \sim \text{EXT}(z_q \mathbf{x}_q)$  and  $f_{t_q}(\mathbf{x}) = \mathbb{E}_{S \sim \text{EXT}(\mathbf{x})}[g_{t_q}(S)]$ . If  $z_q \geq \frac{1}{2}$  and  $\mathbf{u}_q \neq 0$ ,

$$\begin{aligned} \mathbb{E} \left[ \tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q, t_q, z_q, \mathbf{u}_q \right] &= 2(1-1/e) \frac{d}{z_q} \mathbb{E}[g_{t_q}(S_{t_q}) \mid \mathcal{H}_{q-1}, \mathbf{v}_q, t_q, z_q, \mathbf{u}_q] \\ &= 2(1-1/e) \frac{d}{z_q} f_{t_q}(z_q \mathbf{x}_q + z_q \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle \mathbf{u}_q). \end{aligned}$$

Condition on  $\mathcal{H}_{q-1}, \mathbf{v}_q, t_q, z_q$ , then

$$\begin{aligned} \mathbb{E} \left[ \tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q, t_q, z_q \right] &= \frac{1}{2} (-2(1-1/e) \frac{d}{z_q} f_{t_q}(z_q \mathbf{x}_q)) + \sum_{i=1}^d \frac{1}{d} 2(1-1/e) \frac{d}{z_q} f_{t_q}(z_q \mathbf{x}_q + z_q \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle \mathbf{e}_i) \\ &= (1-1/e) \langle \mathbf{H}_q \mathbf{v}_q, \nabla f_{t_q}(z_q \cdot \mathbf{x}_q) \rangle \end{aligned}$$

where the last equality is already proved in the proof of [Lemma 3.1](#).

If  $z_q < \frac{1}{2}$ ,

$$\begin{aligned} &\mathbb{E} \left[ \tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q, z_q, t_q, \mathbf{u}_q, \mathbf{y}_{t_q} \right] \\ &= \begin{cases} -4(1-1/e)d \cdot \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle \mathbb{E}[g_{t_q}(S_{t_q}) \mid \mathcal{H}_{q-1}, \mathbf{v}_q, z_q, t_q, \mathbf{u}_q, \mathbf{y}_{t_q}] & \text{if } \mathbf{y}_{t_q} = z_q \mathbf{x}_q, \\ 4(1-1/e)d \cdot \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle \mathbb{E}[g_{t_q}(S_{t_q}) \mid \mathcal{H}_{q-1}, \mathbf{v}_q, z_q, t_q, \mathbf{u}_q, \mathbf{y}_{t_q}] & \text{if } \mathbf{y}_{t_q} = z_q \mathbf{x}_q + \frac{1}{2} \mathbf{u}_q. \end{cases} \\ &= \begin{cases} -4(1-1/e)d \cdot \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle f_{t_q}(\mathbf{y}_{t_q}) & \text{if } \mathbf{y}_{t_q} = z_q \mathbf{x}_q, \\ 4(1-1/e)d \cdot \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle f_{t_q}(\mathbf{y}_{t_q}) & \text{if } \mathbf{y}_{t_q} = z_q \mathbf{x}_q + \frac{1}{2} \mathbf{u}_q. \end{cases} \end{aligned}$$

Condition on  $\mathcal{H}_{q-1}, \mathbf{v}_q, t_q, z_q$ ,

$$\begin{aligned} &\mathbb{E} \left[ \tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q, t_q, z_q \right] \\ &= \sum_{i=1}^d \frac{1}{d} \left( \frac{1}{2} \cdot (-4(1-1/e)d \cdot \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle f(z_q \mathbf{x}_q)) + \frac{1}{2} \cdot (4(1-1/e)d \cdot \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle f(z_q \mathbf{x}_q + \frac{1}{2} \mathbf{e}_i)) \right) \\ &= \sum_{i=1}^d 2(1-1/e) \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle \left( f(z_q \mathbf{x}_q + \frac{1}{2} \mathbf{e}_i) - f(z_q \mathbf{x}_q) \right) \end{aligned}$$



$$\begin{aligned}
 &= \sum_{i=1}^d 2(1-1/e) \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle \frac{1}{2} \frac{\partial f}{\partial x_i}(z_q \mathbf{x}_q) \\
 &= (1-1/e) \sum_{i=1}^d \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{e}_i \rangle \langle \mathbf{e}_i, \nabla f(z_q \mathbf{x}_q) \rangle \\
 &= (1-1/e) \langle \mathbf{H}_q \mathbf{v}_q, \nabla f_{t_q}(z_q \cdot \mathbf{x}_q) \rangle
 \end{aligned}$$

Combining with the case  $z_q \geq \frac{1}{2}$ , we proved this equation whatever  $z_q$  is:

$$\mathbb{E} \left[ \tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q, t_q, z_q \right] = (1-1/e) \langle \mathbf{H}_q \mathbf{v}_q, \nabla f_{t_q}(z_q \cdot \mathbf{x}_q) \rangle$$

Then follow the calculation in [Lemma 3.1](#), we can prove

$$\mathbb{E}[\tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q] = l_q(\mathbf{H}_q \mathbf{v}_q).$$

□

**Lemma 5.1.** For a finite set  $\mathcal{S}$ , and a function family  $\mathcal{G} \subseteq \mathcal{S}^{\mathbb{R}^+}$ , where  $\mathcal{S}^{\mathbb{R}^+}$  is the set of all functions that map element in  $\mathcal{S}$  to  $\mathbb{R}^+$ . If there is an extension mapping  $EXT: \mathcal{K} \rightarrow \Delta(\mathcal{S})$  satisfying following conditions:

1.  $\mathcal{K} \subseteq \mathbb{R}^d$  is a product of standard simplexes.
2. For any  $g \in \mathcal{G}$ ,  $f(\mathbf{x}) = \mathbb{E}_{S \in EXT(\mathbf{x})}[g(S)]$  is a multi-linear, monotone, DR-submodular function, and  $f$  is  $L_1$ -lipschitz continuous,  $f(\mathbf{0}) = 0$ .
3. For any  $s \in \mathcal{S}$ . Here exist  $\mathbf{x} \in \mathcal{K}$  such that  $EXT(\mathbf{x}) = \mathbf{1}_s$ . Where  $\mathbf{1}_s$  assign probability 1 to  $S$  and 0 to other elements of  $\mathcal{S}$ .

then the algorithm `MLSMWrapper` attains expected  $(1-1/e)$ -regret  $\mathcal{R}_{1-1/e}(T) \leq O(d^{5/3} T^{2/3} \log(T))$  on  $(\mathcal{S}, \mathcal{G})$ -bandit.

*Proof of Lemma 5.1.* We first note that  $g_t(S_t)$  is an unbiased estimator of  $f_t(\mathbf{y}_t)$  by the definition of  $f_t$ . The analysis is the same as the analysis of [Algorithm 1](#) except that [Algorithm 2](#) is actually using a new estimator  $\tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q)$  for  $l(\mathbf{H}_q \mathbf{v}_q)$ . We first bound this new estimator.

If  $z_q \geq \frac{1}{2}$ ,

$$\begin{aligned}
 |\tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q)| &= 2(1-1/e) d \frac{g_{t_q}(S_{t_q})}{z_q} \\
 &\leq 4(1-1/e) dM
 \end{aligned}$$

If  $z_q < \frac{1}{2}$ ,

$$\begin{aligned}
 |\tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q)| &= 4(1-1/e) d \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle g_{t_q}(S_{t_q}) \\
 &\leq 4(1-1/e) dM
 \end{aligned}$$

The inequality is because that  $z_q \mathbf{x}_q + \langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle \mathbf{u}_q \in \mathcal{K}$ , and  $\mathbf{u}_q$  is a basis vector. Therefore  $\langle \mathbf{H}_q \mathbf{v}_q, \mathbf{u}_q \rangle \leq D_\infty$ , here the  $D_\infty$  is the  $\infty$ -norm diameter of  $\mathcal{K}$ . Since  $\mathcal{K}$  is a cartesian product of standard simplexes,  $D_\infty = 1$ .

Let  $\tilde{\nabla} \bar{F}'_q(\mathbf{x}_q)$  be the estimator replacing  $\tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q)$  with  $\tilde{l}_q(\mathbf{H}_q \mathbf{v}_q)$ , that is

$$\tilde{\nabla} \bar{F}'_q(\mathbf{x}_q) = d \cdot \tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q) \mathbf{H}_q^{-1} \mathbf{v}_q$$

We bound the dual local norm of  $\tilde{\nabla} \bar{F}'_q(\mathbf{x}_q)$

$$\mathbb{E} \left[ \|\tilde{\nabla} \bar{F}'_q(\mathbf{x}_q)\|_{\Phi, \mathbf{x}_q, * } \mid \mathcal{H}_{q-1} \right] = \mathbb{E} \left[ d^2 \cdot (\tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q))^2 \mathbf{v}_s^T \mathbf{H}_q^{-1} \Phi(\mathbf{x}_q) \mathbf{H}_q^{-1} \mathbf{v}_q \mid \mathcal{H}_{q-1} \right]$$

$$\begin{aligned}
 &\leq 16(1-1/e)^2 d^4 M^2 \mathbb{E} \left[ \mathbf{v}_q^T \Phi(\mathbf{x}_q)^{-1/2} \Phi(\mathbf{x}_q) \Phi(\mathbf{x}_q)^{-1/2} \mathbf{v}_q \mid \mathcal{H}_{q-1} \right] \\
 &\leq 16(1-1/e)^2 d^4 M^2 \|\mathbf{v}_q\|^2 \\
 &\leq 16(1-1/e)^2 d^4 M^2
 \end{aligned}$$

Since we have proved  $\mathbb{E}[\tilde{l}'_q(\mathbf{H}_q \mathbf{v}_q) \mid \mathcal{H}_{q-1}, \mathbf{v}_q] = l_q(\mathbf{H}_q \mathbf{v}_q)$ , then follow the proof of [Lemma 3.2](#) (i), we can prove

$$\mathbb{E} \left[ \tilde{\nabla} \bar{F}'_q(\mathbf{x}_q) \mid \mathcal{H}_{q-1} \right] = \nabla \bar{F}_q(\mathbf{x}_q)$$

Then follow the proof of [Theorem 3.3](#), we have for any  $\mathbf{x}^* \in \mathcal{K}$ ,

$$\begin{aligned}
 \mathbb{E} \left[ \sum_{t=1}^T (1-1/e) f_t(\mathbf{x}^*) - f_t(\mathbf{y}_t) \right] &\leq \eta L \sum_{q=1}^Q \mathbb{E} \left[ \|\tilde{\nabla} \bar{F}'_q(\mathbf{x}_q)\|_{\Phi, \mathbf{x}_q}^2 \mid \mathcal{H}_{q-1} \right] + (1-1/e) L_1 D + MQ + \frac{\nu L \log(\frac{1}{\delta})}{\eta} \\
 &\leq 16(1-1/e)^2 M^2 d^4 \eta T + (1-1/e) L_1 D + MQ + \frac{\nu L \log(T)}{\eta}
 \end{aligned}$$

Let  $S^* = \operatorname{argmax}_{S \in \mathcal{S}} \sum_{t=1}^T g_t(S)$ ,  $\mathbf{x}^*$  be the point satisfies  $\operatorname{EXT}(\mathbf{x}^*) = \mathbf{1}_{S^*}$ , then  $f_t(\mathbf{x}^*) = g_t(S^*)$ .

Let  $\mathcal{H}'_t$  be the history of the first  $t$ -rounds, including the realization of  $\mathbf{v}_q, t_q, z_q, \mathbf{u}_q, \forall q \leq \lceil \frac{t}{L} \rceil$  and the realization of  $S_k, \forall k \leq t$ . Since  $f_t(\mathbf{y}_t) = \mathbb{E}[g_t(S_t) \mid \mathcal{H}_{t-1}]$ , we have,

$$\begin{aligned}
 \mathbb{E} \left[ \sum_{t=1}^T (1-1/e) g_t(S^*) - g_t(S_t) \right] &= \mathbb{E} \left[ \sum_{t=1}^T (1-1/e) f_t(\mathbf{x}^*) - \mathbb{E}[g_t(S_t) \mid \mathcal{H}_{t-1}] \right] \\
 &= \mathbb{E} \left[ \sum_{t=1}^T (1-1/e) f_t(\mathbf{x}^*) - f_t(\mathbf{y}_t) \right] \\
 &\leq 16(1-1/e)^2 M^2 d^4 \eta T + (1-1/e) L_1 + MQ + \frac{\nu L \log(T)}{\eta}
 \end{aligned}$$

If we use the self-concordant barrier described in [Appendix D](#) as the input  $\Phi$  here, by [Lemma D.1](#),  $\nu = O(d)$ . Then we set  $\eta = d^{-7/3} T^{-1/3}$ ,  $L = d^{-5/3} T^{1/3}$ ,  $Q = T/L = d^{5/3} T^{2/3}$ . Then

$$\mathcal{R}_{1-1/e}(T) = \mathbb{E} \left[ \sum_{t=1}^T (1-1/e) g_t(S^*) - g_t(S_t) \right] = O\left(d^{5/3} T^{2/3} \log(T)\right).$$

□

## E.2. Proof of [Lemma 5.3](#) and [Corollary 5.4](#)

Before proving [Lemma 5.3](#), we first prove a useful lemma.

**Lemma E.2.** *Let  $g : 2^G \rightarrow \mathbb{R}_+$  be a monotone submodular set function,  $S$  is a subset of  $G$ ,  $s_1, s_2 \in G$  and there is no any other restriction on  $s_1$  and  $s_2$ , they may be the same element or not, and they may be in  $S$  or not. Then  $g(S \cup \{s_1, s_2\}) - g(S \cup \{s_1\}) \leq g(S \cup \{s_2\}) - g(S)$ .*

*Proof.* If  $s_1 \in S$  and  $s_2 \notin S$ , then  $g(S \cup \{s_1, s_2\}) - g(S \cup \{s_1\}) = g(S \cup \{s_2\}) - g(S)$ . If  $s_1 \notin S$  and  $s_2 \in S$ , then  $g(S \cup \{s_1, s_2\}) - g(S \cup \{s_1\}) = g(S \cup \{s_1\}) - g(S \cup \{s_1\}) = 0$  and  $g(S \cup \{s_2\}) - g(S) = g(S) - g(S) = 0$ . If  $s_1 \in S, s_2 \in S$ , then  $g(S \cup \{s_1, s_2\}) - g(S \cup \{s_1\}) = g(S \cup \{s_2\}) - g(S) = g(S) - g(S) = 0$ , the inequality holds.

If  $s_1 \notin S, s_2 \notin S$  and  $s_1 \neq s_2$ , then the result holds due to the submodularity of  $g$ . If  $s_1 = s_2 \notin S$ , then  $g(S \cup \{s_1, s_2\}) - g(S \cup \{s_1\}) = 0$  and  $g(S \cup \{s_2\}) - g(S) \geq 0$  by the monotonicity of  $g$ . □

**Lemma 5.3.** *For  $\mathcal{G}_{MS}$ , the extension mapping  $\operatorname{EXT}_{PM} : \mathcal{K} \rightarrow \Delta(\mathcal{S}_{PM})$  satisfies the conditions in [Lemma 5.1](#). Moreover,  $\mathcal{K}$  is in a  $\sum_{k=1}^K r_k |G_k|$  dimensional real vector space. For any  $g \in \mathcal{G}_{MS}$ , the continuous extension  $f(\mathbf{x}) = \mathbb{E}_{s \in \operatorname{EXT}(\mathbf{x})}[g(s)]$  is  $M \sqrt{\sum_{k=1}^K r_k |G_k|}$ -lipschitz.*

*Proof of Lemma 5.3.* We first prove that for any  $S \in \mathcal{S}_{PM}$ , there is a  $\mathbf{x} \in \mathcal{K}$  such that  $\text{EXT}_{PM}(\mathbf{x}) = \mathbf{1}_S$ . For any  $S \in \mathcal{S}_{PM}$ , it can be partitioned into  $S = \cup_{k=1}^K S_k$  such that  $S_k \subseteq G_k$  and  $|S_k| \leq r_k$ . For any  $k$ , we select  $|S_k|$  standard simplexes  $\Delta_{G_k}^{k,i}$ ,  $i \in [|S_k|]$ , and we assign probability 1 to the elements in  $S_k$  respectively in these standard simplexes. Thus the condition 1 and 3 of Lemma 5.1 is satisfied, the dimension of  $\mathcal{K}$  is obvious. It's enough to show that  $f(\mathbf{x})$  satisfies the condition 2.

$\text{EXT}_{PM}(\mathbf{0})$  assigns probability 1 to the empty set, thus  $f(\mathbf{0}) = 0$ .

Now we check the multi-linearity of  $f$ . Consider a sample  $\omega \in \Omega$ , the probability of  $\omega$  is  $\Pr(\omega) = \prod_{k=1}^K \prod_{i=1}^{r_k} \Pr(\omega_{k,i})$ , and  $\Pr(\omega_{k,i}) \in \{x_{k,i,s} \mid s \in G_k\} \cup \{1 - \sum_{s \in G_k} x_{k,i,s}\}$ , thus  $\Pr(\omega)$  is multi-linear with respect to the variables  $x_{k,i,s}$ . Then we write  $f(\mathbf{x})$  as follows,

$$\begin{aligned} f(\mathbf{x}) &= \mathbb{E}_{S \sim \text{EXT}_{PM}(\mathbf{x})} [g(S)] \\ &= \sum_{S \in \mathcal{S}_{PM}} \Pr(S) g(S) \\ &= \sum_{S \in \mathcal{S}_{PM}} \sum_{\omega \in \rho^{-1}(S)} \Pr(\omega) g(S) \\ &= \sum_{\omega \in \Omega} \Pr(\omega) g(\rho(\omega)) \end{aligned}$$

Since  $g(\rho(\omega))$  is a constant independent from  $\mathbf{x}$ ,  $f(\mathbf{x})$  is a linear combination of multi-linear functions, thus  $f(\mathbf{x})$  is also multi-linear.

Since  $f(\mathbf{x})$  is multi-linear, its partial derivative is

$$\frac{\partial f}{\partial x_{k,i,s}}(\mathbf{x}) = \frac{f(\mathbf{x} \vee (1 - \sum_{s' \in G_k, s' \neq s} x_{k,i,s'}) \mathbf{e}_{k,i,s}) - f(\mathbf{x} \wedge \bar{\mathbf{e}}_{k,i,s})}{1 - \sum_{s' \in G_k, s' \neq s} x_{k,i,s'}}$$

Here  $\wedge$  is the coordinate-wise minimal, and  $\vee$  is the coordinate-wise maximal.  $\mathbf{e}_{k,i,s}$  is the basis vector which takes 1 only for the component indexed  $(k, i, s)$ , and 0 for the other components.  $\bar{\mathbf{e}}_{k,i,s}$  is the vector that take 0 for the component indexed  $(k, i, s)$  and 1 for the other components.

We define two mappings  $\rho_{\vee}^{k,i,s}, \rho_{\wedge}^{k,i,s} : \Omega \rightarrow \Omega$ . For  $\omega = (\omega_{k',i'})_{k',i' \in \Gamma}$ ,  $\Gamma = \{(k', i') \mid 1 \leq k' \leq K, 1 \leq i' \leq r_{k'}, i, k \in \mathbb{N}\}$ .  $\rho_{\vee}^{k,i,s}$  and  $\rho_{\wedge}^{k,i,s}$  only change the component indexed  $(k, i)$ , for any  $(k', i') \neq (k, i)$ ,  $(\rho_{\vee}^{k,i,s}(\omega))_{k',i'} = (\rho_{\wedge}^{k,i,s}(\omega))_{k',i'} = \omega_{k',i'}$ . For the component indexed  $(k, i)$ , let

$$\left( \rho_{\vee}^{k,i,s}(\omega) \right)_{k,i} = \begin{cases} s & \text{if } \omega_{k,i} = \circ \\ \omega_{k,i} & \text{if } \omega_{k,i} \neq \circ \end{cases} \quad (21)$$

and

$$\left( \rho_{\wedge}^{k,i,s}(\omega) \right)_{k,i} = \begin{cases} \circ & \text{if } \omega_{k,i} = s \\ \omega_{k,i} & \text{if } \omega_{k,i} \neq s \end{cases} \quad (22)$$

We make the following important claim

**Claim E.3.**

$$f\left(\mathbf{x} \vee \left(1 - \sum_{s' \in G_k, s' \neq s} x_{k,i,s'}\right) \mathbf{e}_{k,i,s}\right) = \mathbb{E}_{\omega \sim \text{pre-EXT}_{PM}(\mathbf{x})} \left[ g\left(\rho\left(\rho_{\vee}^{k,i,s}(\omega)\right)\right) \right] \quad (23)$$

and

$$f(\mathbf{x} \wedge \bar{\mathbf{e}}_{k,i,s}) = \mathbb{E}_{\omega \sim \text{pre-EXT}_{PM}(\mathbf{x})} \left[ g\left(\rho\left(\rho_{\wedge}^{k,i,s}(\omega)\right)\right) \right]. \quad (24)$$

*proof of Claim E.3.* Let  $\mathbf{x}^{\vee,k,i,s} := \mathbf{x} \vee \left(1 - \sum_{s' \in G_k, s' \neq s} x_{k,i,s'}\right) \mathbf{e}_{k,i,s}$ .

$$f(\mathbf{x}^{\vee,k,i,s}) = \mathbb{E}_{\omega \sim \text{pre-EXT}_{PM}(\mathbf{x}^{\vee,k,i,s})} [g(\rho(\omega))]$$

Let  $\omega \sim \text{pre-EXT}_{PM}(\mathbf{x})$ , it's enough to show  $\rho_{\vee}^{k,i,s}(\omega) \sim \text{pre-EXT}(\mathbf{x}^{\vee,k,i,s})$ . Let  $\omega_{\vee} \sim \text{pre-EXT}_{PM}(\mathbf{x}^{\vee,k,i,s})$ . For  $(k', i') \neq (k, i)$  and  $s' \in G_{k'}$ ,  $\Pr((\omega_{\vee})_{k',i'} = s') = \Pr((\rho_{\vee}^{k,i,s}(\omega))_{k',i'} = s') = x_{k',i',s'}$ , and  $\Pr((\omega_{\vee})_{k',i'} = \circ) = \Pr((\rho_{\vee}^{k,i,s}(\omega))_{k',i'} = \circ) = 1 - \sum_{s' \in G_{k'}} x_{k',i',s'}$ . For  $(k, i)$ -component and  $s' \neq s$ ,  $\Pr((\omega_{\vee})_{k,i} = s') = x_{k,i,s'} = \Pr((\rho_{\vee}^{k,i,s}(\omega))_{k,i} = s')$ . For  $s$  and  $\circ$ ,

$$\Pr((\omega_{\vee})_{k,i} = s) = 1 - \sum_{s' \in G_k, s' \neq s} x_{k,i,s'}, \quad \Pr((\omega_{\vee})_{k,i} = \circ) = 0$$

$(\rho_{\vee}^{k,i,s}(\omega))_{k,i} = s$  whenever  $\omega \in \{s, \circ\}$ , thus

$$\Pr((\rho_{\vee}^{k,i,s}(\omega))_{k,i} = s) = 1 - \sum_{s' \in G_k, s' \neq s} x_{k,i,s'} = \Pr((\omega_{\vee})_{k,i} = s)$$

Since  $\rho_{\vee}^{k,i,s}(\omega)_{k,i}$  never be  $\circ$ ,  $\Pr(\rho_{\vee}^{k,i,s}(\omega)_{k,i} = \circ) = 0$ . Thus  $\rho_{\vee}^{k,i,s}(\omega) \sim \text{pre-EXT}(\mathbf{x}^{\vee,k,i,s})$ , and

$$f(\mathbf{x}^{\vee,k,i,s}) = \mathbb{E}_{\omega \sim \text{pre-EXT}_{PM}(\mathbf{x}^{\vee,k,i,s})} [g(\rho(\omega))] = \mathbb{E}_{\omega \sim \text{pre-EXT}_{PM}(\mathbf{x})} \left[ g \left( \rho \left( \rho_{\vee}^{k,i,s}(\omega) \right) \right) \right].$$

In brief, whenever  $\omega_{k,i} \in \{s, \circ\}$ ,  $\rho_{\vee}^{k,i,s}(\omega) = s$ . That is,  $\Pr((\rho_{\vee}^{k,i,s}(\omega))_{k,i} = s) = \Pr(\omega_{k,i} \in \{s, \circ\}) = 1 - \Pr(\omega_{k,i} \notin \{s, \circ\}) = 1 - \sum_{s' \neq s, s' \in G_k} x_{k,i,s'}$  which is the same as a sample in  $\text{pre-EXT}_{PM}(\mathbf{x}^{\vee,k,i,s})$ .

For (24), we can define  $\mathbf{x}^{\wedge,k,i,s} := \mathbf{x} \bar{\mathbf{e}}_{k,i,s}$  and let  $\omega_{\wedge} \sim \text{pre-EXT}_{PM}(\mathbf{x}^{\wedge,k,i,s})$ . One can check that for  $(k', i', s') \neq (k, i, s)$ ,  $\Pr((\rho_{\wedge}^{k,i,s}(\omega))_{k',i'} = s') = \Pr((\omega_{\wedge})_{k',i'} = s') = x_{k',i',s'}$  and  $\Pr((\rho_{\wedge}^{k,i,s}(\omega))_{k',i'} = \circ) = \Pr((\omega_{\wedge})_{k',i'} = \circ) = 1 - \sum_{s' \in G_{k'}} x_{k',i',s'}$ . For  $(k, i, s)$ ,

$$\Pr((\rho_{\wedge}^{k,i,s}(\omega))_{k,i} = \circ) = \Pr((\omega_{\wedge})_{k,i} = \circ) = 1 - \sum_{s' \in G_k, s' \neq s} x_{k,i,s'}$$

and

$$\Pr((\rho_{\wedge}^{k,i,s}(\omega))_{k,i} = s) = \Pr((\omega_{\wedge})_{k,i} = s) = 0.$$

So  $\rho_{\wedge}^{k,i,s}(\omega) \sim \text{pre-EXT}(\mathbf{x}^{\wedge,k,i,s})$  and (24) holds.  $\square$

If  $\omega_{k,i} = \circ$  or  $\omega_{k,i} = s$ , one can check that  $(\rho_{\vee}^{k,i,s}(\omega))_{k,i} = s$  and  $(\rho_{\wedge}^{k,i,s}(\omega))_{k,i} = \omega_{k,i} = \circ$ , therefore  $\rho(\rho_{\vee}^{k,i,s}(\omega)) = \rho(\rho_{\wedge}^{k,i,s}(\omega)) \cup \{s\}$ , so  $\rho(\rho_{\wedge}^{k,i,s}(\omega)) \subseteq \rho(\rho_{\vee}^{k,i,s}(\omega))$ . if  $\omega_{k,i} \notin \{\circ, s\}$ , then  $\rho_{\vee}^{k,i,s}(\omega) = \rho_{\wedge}^{k,i,s}(\omega)$ , so  $\rho(\rho_{\wedge}^{k,i,s}(\omega)) = \rho(\rho_{\vee}^{k,i,s}(\omega))$ . Thus we proved  $\rho(\rho_{\wedge}^{k,i,s}(\omega)) \subseteq \rho(\rho_{\vee}^{k,i,s}(\omega))$  for any  $\omega$ . Since  $g$  is monotone, we have

$$\frac{\partial f}{\partial x_{k,i,s}}(\mathbf{x}) = \frac{\mathbb{E}_{\omega \sim \text{pre-EXT}_{PM}(\mathbf{x})} \left[ g \left( \rho \left( \rho_{\vee}^{k,i,s}(\omega) \right) \right) - g \left( \rho \left( \rho_{\wedge}^{k,i,s}(\omega) \right) \right) \right]}{1 - \sum_{s' \in G_k, s' \neq s} x_{k,i,s'}} \geq 0$$

Note that  $\rho_{\wedge}^{k,i,s}(\omega) \neq \rho_{\vee}^{k,i,s}(\omega)$  only happens when  $\omega_{k,i} = s$  or  $\omega_{k,i} = \circ$ , thus,

$$\begin{aligned} \left| \frac{\partial f}{\partial x_{k,i,s}}(\mathbf{x}) \right| &\leq \frac{M \Pr(\omega_{k,i} \in \{s, \circ\})}{1 - \sum_{s' \in G_k, s' \neq s} x_{k,i,s'}} \\ &= \frac{M(1 - \sum_{s' \in G_k, s' \neq s} x_{k,i,s'})}{1 - \sum_{s' \in G_k, s' \neq s} x_{k,i,s'}} = M \end{aligned}$$

which shows that  $\|\nabla f\|_\infty \leq M$ . Since  $\nabla f \in \mathbb{R}^{\sum_{k=1}^K r_k |G_k|}$ ,  $\|\nabla f\|_2 \leq \sqrt{\sum_{k=1}^K r_k |G_k|} \|\nabla f\|_\infty = M \sqrt{\sum_{k=1}^K r_k |G_k|}$ . Therefore,  $f$  is  $M \sqrt{\sum_{k=1}^K r_k |G_k|}$ -lipschitz continuous.

Since the partial derivative of a multi-linear function is also multi-linear,  $\frac{\partial f}{\partial x_{k,i,s}}(\mathbf{x})$  is multi-linear for any  $k, i, s$ . Then the second derivative of  $f$  can be written as

$$\begin{aligned} \frac{\partial^2 f}{\partial x_{k_1, i_1, s_1} \partial x_{k_2, i_2, s_2}}(\mathbf{x}) &= \frac{\frac{\partial f}{\partial x_{k_2, i_2, s_2}}(\mathbf{x} \vee (1 - \sum_{s' \in G_{k_1}, s' \neq s_1} x_{k_1, i_1, s'}) \mathbf{e}_{k_1, i_1, s_1}) - \frac{\partial f}{\partial x_{k_2, i_2, s_2}}(\mathbf{x} \wedge \bar{\mathbf{e}}_{k_1, i_1, s_1})}{1 - \sum_{s' \in G_{k_1}, s' \neq s_1} x_{k_1, i_1, s'}} \\ &= \frac{f(\mathbf{x} \vee (1 - \sum_{s' \in G_{k_1}, s' \neq s_1} x_{k_1, i_1, s'}) \mathbf{e}_{k_1, i_1, s_1} \vee (1 - \sum_{s' \in G_{k_2}, s' \neq s_2} x_{k_2, i_2, s'}) \mathbf{e}_{k_2, i_2, s_2})}{(1 - \sum_{s' \in G_{k_1}, s' \neq s_1} x_{k_1, i_1, s'}) (1 - \sum_{s' \in G_{k_2}, s' \neq s_2} x_{k_2, i_2, s'})} \\ &\quad - \frac{f((\mathbf{x} \wedge \bar{\mathbf{e}}_{k_2, i_2, s_2}) \vee (1 - \sum_{s' \in G_{k_1}, s' \neq s_2} x_{k_1, i_1, s'}) \mathbf{e}_{k_1, i_1, s_1})}{(1 - \sum_{s' \in G_{k_1}, s' \neq s_1} x_{k_1, i_1, s'}) (1 - \sum_{s' \in G_{k_2}, s' \neq s_2} x_{k_2, i_2, s'})} \\ &\quad - \frac{f((\mathbf{x} \vee (1 - \sum_{s' \in G_{k_2}, s' \neq s_2} x_{k_2, i_2, s'}) \mathbf{e}_{k_2, i_2, s_2}) \wedge \bar{\mathbf{e}}_{k_1, i_1, s_1})}{(1 - \sum_{s' \in G_{k_2}, s' \neq s_2} x_{k_2, i_2, s'}) (1 - \sum_{s' \in G_{k_1}, s' \neq s_1} x_{k_1, i_1, s'})} \\ &\quad + \frac{f(\mathbf{x} \wedge \bar{\mathbf{e}}_{k_2, i_2, s_2} \wedge \bar{\mathbf{e}}_{k_1, i_1, s_1})}{(1 - \sum_{s' \in G_{k_2}, s' \neq s_2} x_{k_2, i_2, s'}) (1 - \sum_{s' \in G_{k_1}, s' \neq s_1} x_{k_1, i_1, s'})} \end{aligned}$$

To prove the DR-submodularity of  $f$ , We define 4 mappings  $\rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}$ ,  $\rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}$ ,  $\rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}$ ,  $\rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2} : \Omega \rightarrow \Omega$ .

$$\begin{aligned} \rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) &= \rho_{\vee}^{k_1, i_1, s_1}(\rho_{\vee}^{k_2, i_2, s_2}(\omega)) & \rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) &= \rho_{\vee}^{k_1, i_1, s_1}(\rho_{\wedge}^{k_2, i_2, s_2}(\omega)) \\ \rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) &= \rho_{\wedge}^{k_1, i_1, s_1}(\rho_{\vee}^{k_2, i_2, s_2}(\omega)) & \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) &= \rho_{\wedge}^{k_1, i_1, s_1}(\rho_{\wedge}^{k_2, i_2, s_2}(\omega)) \end{aligned}$$

Same as Claim E.3, we have

$$\begin{aligned} &\frac{\partial f}{\partial x_{k_1, i_1, s_1} \partial x_{k_2, i_2, s_2}}(\mathbf{x}) \\ &= \frac{\mathbb{E}_{\omega \sim \text{pre-EXT}_{PM}(\mathbf{x})} \left[ \rho(\rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)) - \rho(\rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)) \right]}{(1 - \sum_{s' \in G_{k_2}, s' \neq s_2} x_{k_2, i_2, s'}) (1 - \sum_{s' \in G_{k_1}, s' \neq s_1} x_{k_1, i_1, s'})} \\ &\quad + \frac{\mathbb{E}_{\omega \sim \text{pre-EXT}_{PM}(\mathbf{x})} \left[ -\rho(\rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)) + \rho(\rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)) \right]}{(1 - \sum_{s' \in G_{k_2}, s' \neq s_2} x_{k_2, i_2, s'}) (1 - \sum_{s' \in G_{k_1}, s' \neq s_1} x_{k_1, i_1, s'})}. \end{aligned}$$

We first consider the situation where  $k_1 = k_2$  and  $i_1 = i_2$ . Recall that  $\Pr(\omega_{k_1, i_1}) \in \{x_{k_1, i_1, s} \mid s \in G_k\} \cup \{1 - \sum_{s \in G_{k_1}} x_{k_1, i_1, s}\}$ , so  $\frac{\partial f}{\partial x_{k_1, i_1, s_1} \partial x_{k_2, i_2, s_2}} \Pr(\omega_{k_1, i_1}) = 0$  when  $k_1 = k_2$  and  $i_1 = i_2$ . In this situation,

$$\begin{aligned} \frac{\partial f}{\partial x_{k_1, i_1, s_1} \partial x_{k_2, i_2, s_2}}(\mathbf{x}) &= \sum_{\omega \in \Omega} g(\rho(\omega)) \frac{\partial f}{\partial x_{k_1, i_1, s_1} \partial x_{k_2, i_2, s_2}} \Pr(\omega) \\ &= \sum_{\omega \in \Omega} g(\rho(\omega)) \left( \prod_{\substack{k' \in [K], i' \in [r_{k'}] \\ k' \neq k_1 \text{ or } i' \neq i_1}} \Pr(\omega_{k', i'}) \right) \frac{\partial f}{\partial x_{k_1, i_1, s_1} \partial x_{k_2, i_2, s_2}} \Pr(\omega_{k_1, i_1}) \\ &= 0 \quad \text{if } k_1 = k_2 \text{ and } i_1 = i_2 \end{aligned} \tag{25}$$

Then we consider the situation that  $k_1 \neq k_2$  or  $i_1 \neq i_2$ ,

**Case 1 :**  $k_1 \neq k_2$  or  $i_1 \neq i_2$ ,  $\omega_{k_1, i_1} \in \{s_1, \circ\}$  and  $\omega_{k_2, i_2} \in \{s_2, \circ\}$ . In this case,

$$\begin{aligned} \left( \rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right)_{k_1, i_1} &= s_1 & \left( \rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right)_{k_2, i_2} &= s_2 \\ \left( \rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right)_{k_1, i_1} &= s_1 & \left( \rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right)_{k_2, i_2} &= \circ \\ \left( \rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right)_{k_1, i_1} &= \circ & \left( \rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right)_{k_2, i_2} &= s_2 \\ \left( \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right)_{k_1, i_1} &= \circ & \left( \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right)_{k_2, i_2} &= \circ \end{aligned}$$

$\rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)$ ,  $\rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)$ ,  $\rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)$ ,  $\rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)$  are equal in all but above two components  $k_1, i_1$  and  $k_2, i_2$ . Thus

$$\begin{aligned} \rho \left( \rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) &= \rho \left( \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) \cup \{s_1, s_2\} \\ \rho \left( \rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) &= \rho \left( \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) \cup \{s_1\} \\ \rho \left( \rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) &= \rho \left( \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) \cup \{s_2\} \end{aligned}$$

For monotone submodular  $g$ , by Lemma E.2,

$$\begin{aligned} g \left( \rho \left( \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) \cup \{s_1, s_2\} \right) &- g \left( \rho \left( \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) \cup \{s_1\} \right) \\ &\geq g \left( \rho \left( \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) \cup \{s_2\} \right) - g \left( \rho \left( \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) \right) \end{aligned}$$

thus,

$$\rho \left( \rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) - \rho \left( \rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) - \rho \left( \rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) + \left( \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) \leq 0$$

**Case 2 :**  $k_1 \neq k_2$  or  $i_1 \neq i_2$ ,  $\omega_{k_1, i_1} \in \{s_1, \circ\}$  and  $\omega_{k_2, i_2} \notin \{s_2, \circ\}$ . In this case,

$$\rho_{\vee}^{k_2, i_2, s_2}(\omega) = \rho_{\wedge}^{k_2, i_2, s_2}(\omega) = \omega$$

thus,

$$\begin{aligned} \rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) &= \rho_{\vee}^{k_1, i_1, s_1}(\omega) & \rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) &= \rho_{\vee}^{k_1, i_1, s_1}(\omega) \\ \rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) &= \rho_{\wedge}^{k_1, i_1, s_1}(\omega) & \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) &= \rho_{\wedge}^{k_1, i_1, s_1}(\omega) \end{aligned}$$

and

$$\rho \left( \rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) - \rho \left( \rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) - \rho \left( \rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) + \left( \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) = 0$$

**Case 3 :**  $k_1 \neq k_2$  or  $i_1 \neq i_2$ ,  $\omega_{k_1, i_1} \notin \{s_1, \circ\}$  and  $\omega_{k_2, i_2} \in \{s_2, \circ\}$ . Since  $(k_1, i_1) \neq (k_2, i_2)$ ,  $\rho_{\vee}^{k_2, i_2, s_2}$  and  $\rho_{\wedge}^{k_2, i_2, s_2}$  do not change the  $k_1, i_1$  component of  $\omega$ , that is,

$$\left( \rho_{\vee}^{k_2, i_2, s_2}(\omega) \right)_{k_1, i_1} = \left( \rho_{\wedge}^{k_2, i_2, s_2}(\omega) \right)_{k_1, i_1} = \omega_{k_1, i_1}$$

Therefore,

$$\begin{aligned} \rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) &= \rho_{\vee}^{k_2, i_2, s_2}(\omega) & \rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) &= \rho_{\wedge}^{k_2, i_2, s_2}(\omega) \\ \rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) &= \rho_{\vee}^{k_2, i_2, s_2}(\omega) & \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) &= \rho_{\wedge}^{k_2, i_2, s_2}(\omega) \end{aligned}$$

and

$$\rho \left( \rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) - \rho \left( \rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) - \rho \left( \rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) + \left( \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) \right) = 0$$

**Case 4 :**  $k_1 \neq k_2$  or  $i_1 \neq i_2$ ,  $\omega_{k_1, i_1} \notin \{s_1, \circ\}$  and  $\omega_{k_2, i_2} \notin \{s_2, \circ\}$ . In this case,

$$\rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) = \rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) = \rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) = \rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega) = \omega$$

thus

$$\rho\left(\rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)\right) - \rho\left(\rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)\right) - \rho\left(\rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)\right) + \left(\rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)\right) = 0.$$

In all 4 cases above, whatever  $\omega$  is,

$$\rho\left(\rho_{\vee, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)\right) - \rho\left(\rho_{\vee, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)\right) - \rho\left(\rho_{\wedge, \vee}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)\right) + \left(\rho_{\wedge, \wedge}^{k_1, i_1, s_1, k_2, i_2, s_2}(\omega)\right) \leq 0$$

holds. Thus,

$$\frac{\partial f}{\partial x_{k_1, i_1, s_1} \partial x_{k_2, i_2, s_2}}(\mathbf{x}) \leq 0 \quad \text{if } k_1 \neq k_2 \text{ or } i_1 \neq i_2. \quad (26)$$

Combining (25) and (26),  $\frac{\partial f}{\partial x_{k_1, i_1, s_1} \partial x_{k_2, i_2, s_2}}(\mathbf{x}) \leq 0$  for any  $k_1, i_1, s_1, k_2, i_2, s_2$ , which shows the DR-submodularity of  $f$ .  $\square$

**Corollary 5.4.** *There is an algorithm attaining the expected  $(1 - 1/e)$ -regret of*

$$\mathcal{R}_{1-1/e}(T) \leq O\left(\left(\sum_{k=1}^K r_k |G_k|\right)^{5/3} T^{2/3} \log T\right)$$

on any  $(\mathcal{S}_{PM}, \mathcal{G}_{MS})$ -bandit.

*Proof of Corollary 5.4.* Since  $\text{EXT}_{PM}$  satisfies the conditions in Lemma 5.1 and the dimension of  $\mathcal{K}$  is  $d = \sum_{k=1}^K r_k |G_k|$ . This is a direct corollary of Lemma 5.1.  $\square$

### E.3. Proof of Lemma 5.5 and Corollary 5.6

**Lemma 5.5.** *For  $\mathcal{G}_{SS}$ , the extension mapping  $\text{EXT}_{SS} : \mathcal{K} \rightarrow \Delta(\mathcal{S}_{OL})$  satisfies the conditions in Lemma 5.1. Moreover,  $\mathcal{K}$  is in a  $|G|^2 - |G|$  dimensional real vector space. For any  $g \in \mathcal{G}_{SS}$ , the continuous extension  $f(\mathbf{x}) = \mathbb{E}_{s \in \text{EXT}(\mathbf{x})}[g(s)]$  is  $M|G|$ -lipschitz.*

*Proof of Lemma 5.5.* The condition 1 and 3 of Lemma 5.1 are obviously satisfied. Now we check the multi-linearity of  $f$ . For  $\mathbf{x} \in \mathcal{K}$ , we write  $\mathbf{x} = (x_{i,s})_{i \in |G|, s \in G'}$ . Consider  $S \in \mathcal{S}_{OL}$ , given  $\mathbf{x}$ , the probability of  $S$  in distribution  $\text{EXT}_{SS}(\mathbf{x})$  is  $\Pr(S) = \prod_{i=1}^{|G|} \Pr(S_i)$ . For any  $1 \leq i \leq |G|$ ,  $\Pr(S_i) \in \{x_{i,s} \mid s \in G'\} \cup \{1 - \sum_{s \in G'} x_{i,s}\}$ , thus  $\Pr(S)$  is multi-linear with respect to the variables  $x_{i,s}$ . Then we write  $f(\mathbf{x})$  as follows,

$$\begin{aligned} f(\mathbf{x}) &= \mathbb{E}_{S \sim \text{EXT}_{SS}(\mathbf{x})}[g(S)] \\ &= \sum_{S \in \mathcal{S}_{OL}} \Pr(S) g(S) \end{aligned}$$

Since  $g(S)$  is a constant independent from  $\mathbf{x}$ ,  $f(\mathbf{x})$  is a linear combination of multi-linear functions, thus  $f(\mathbf{x})$  is also multi-linear.

$\text{EXT}(\mathbf{0})$  assigns probability 1 to the ordered list  $\{\circ\}^{|G|}$ , thus  $f(\mathbf{0}) = g(\{\circ\}^{|G|}) = 0$ . Next we check the monotonicity and DR-submodularity of  $f(\mathbf{x})$ .

Define  $\rho_{\vee}^{i,s}, \rho_{\wedge}^{i,s} : \mathcal{S}_{OL} \rightarrow \mathcal{S}_{OL}$ .  $\rho_{\vee}^{i,s}(S) \neq S$  only when the  $i$ -th position of  $S$  is  $\circ$ ,  $\rho_{\vee}^{i,s}(S)$  change the  $i$ -th position of  $S$  to  $s$  and keep other positions unchanged.  $\rho_{\wedge}^{i,s}(S) \neq S$  only when the  $i$ -th position of  $S$  is  $s$ ,  $\rho_{\wedge}^{i,s}(S)$  change the  $i$ -th position of  $S$  to  $\circ$  and keep other positions unchanged. Then,

$$\frac{\partial f}{\partial x_{i,s}}(\mathbf{x}) = \frac{\mathbb{E}_{S \sim \text{EXT}_{SS}(\mathbf{x})} \left[ g(\rho_{\vee}^{i,s}(S)) - g(\rho_{\wedge}^{i,s}(S)) \right]}{1 - \sum_{s' \in G, s' \neq s} x_{i,s'}}$$

Let  $S^{\leq k}$  be the set containing the first  $k$  elements in the ordered list  $S$ . Then  $(\rho_{\wedge}^{i,s}(S))^{\leq k} \subseteq (\rho_{\vee}^{i,s}(S))^{\leq k}$ ,  $\forall k$ , recall  $g_k$  is monotone and  $\lambda_k > 0$  for any  $k$ ,

$$\frac{\partial f}{\partial x_{i,s}}(\mathbf{x}) = \frac{\mathbb{E}_{S \sim \text{EXT}_{SSM}(\mathbf{x})} \left[ \sum_{k=1}^{|G_k|} \lambda_k \left( g_k((\rho_{\vee}^{i,s}(S))^{\leq k}) - g_k((\rho_{\wedge}^{i,s}(S))^{\leq k}) \right) \right]}{1 - \sum_{s' \in G, s' \neq s} x_{i,s'}} \geq 0$$

Thus  $f$  is monotone. Since  $\Pr(\rho_{\vee}^{i,s}(S) \neq \rho_{\wedge}^{i,s}(S)) \leq 1 - \sum_{s' \in G, s' \neq s} x_{i,s'}$  and  $\frac{\partial f}{\partial x_{i,s}}(\mathbf{x}) \leq M$ . Thus  $\|\nabla f(\mathbf{x})\|_{\infty} \leq M$ ,  $\|\nabla f(\mathbf{x})\|_2 \leq M\sqrt{|G|(|G|-1)} \leq M|G|$ ,  $f(\mathbf{x})$  is  $M|G|$ -lipschitz.

We then define

$$\begin{aligned} \rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S) &= \rho_{\vee}^{i_1, s_1} \left( \rho_{\vee}^{i_2, s_2}(S) \right) & \rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S) &= \rho_{\vee}^{i_1, s_1} \left( \rho_{\wedge}^{i_2, s_2}(S) \right) \\ \rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S) &= \rho_{\wedge}^{i_1, s_1} \left( \rho_{\vee}^{i_2, s_2}(S) \right) & \rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S) &= \rho_{\wedge}^{i_1, s_1} \left( \rho_{\wedge}^{i_2, s_2}(S) \right) \end{aligned}$$

Then,

$$\frac{\partial f}{\partial x_{i_1, s_1} \partial x_{i_2, s_2}}(\mathbf{x}) = \frac{\mathbb{E}_{S \sim \text{EXT}_{SSM}(\mathbf{x})} \left[ g(\rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S)) - g(\rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S)) - g(\rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S)) + g(\rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S)) \right]}{(1 - \sum_{s' \in G, s' \neq s_1} x_{i_1, s'}) (1 - \sum_{s' \in G, s' \neq s_2} x_{i_2, s'})}$$

We then prove  $g(\rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S)) - g(\rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S)) - g(\rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S)) + g(\rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S)) \leq 0$  for any  $S \in \mathcal{S}$ . It's enough to prove  $g_i((\rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S))^{\leq i}) - g_i((\rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S))^{\leq i}) - g_i((\rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S))^{\leq i}) + g_i((\rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S))^{\leq i}) \leq 0$  for any  $i \in [|G|]$ . Note that if  $\max\{i_1, i_2\} > i$  then  $g_i((\rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S))^{\leq i}) - g_i((\rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S))^{\leq i}) - g_i((\rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S))^{\leq i}) + g_i((\rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S))^{\leq i}) = 0$ , so we now consider the case  $\max\{i_1, i_2\} \leq i$ .

**Case 1 :**  $i_1 = i_2$ . In this case, since  $\Pr(S_{i_1}) \in \{x_{i_1, s} \mid s \in G'\} \cup \{1 - \sum_{s \in G'} x_{i_1, s}\}$ ,  $\frac{\partial^2 \Pr(S_{i_1})}{\partial x_{i_1, s_1} \partial x_{i_2, s_2}} = 0$  when  $i_1 = i_2$ . we have

$$\begin{aligned} \frac{\partial^2 f}{\partial x_{i_1, s_1} \partial x_{i_2, s_2}}(\mathbf{x}) &= \frac{\partial^2}{\partial x_{i_1, s_1} \partial x_{i_2, s_2}} \sum_{S \in \mathcal{S}_{OL}} g(S) \prod_{i=1}^{|G|} \Pr(S_i) \\ &= \sum_{S \in \mathcal{S}_{OL}} g(S) \left( \prod_{i \neq i_1} \Pr(S_i) \right) \frac{\partial^2 \Pr(S_{i_1})}{\partial x_{i_1, s_1} \partial x_{i_2, s_2}} \\ &= 0 \end{aligned}$$

**Case 2 :**  $i_1 \neq i_2$ ,  $S_{i_1} \in \{s_1, \circ\}$  and  $S_{i_2} \in \{s_2, \circ\}$ . In this case,

$$\begin{aligned} \left( \rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S) \right)_{i_1} &= s_1 & \left( \rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S) \right)_{i_2} &= s_2 \\ \left( \rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S) \right)_{i_1} &= s_1 & \left( \rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S) \right)_{i_2} &= \circ \\ \left( \rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S) \right)_{i_1} &= \circ & \left( \rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S) \right)_{i_2} &= s_2 \\ \left( \rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S) \right)_{i_1} &= \circ & \left( \rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S) \right)_{i_2} &= \circ \end{aligned}$$

$\rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S)$ ,  $\rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S)$ ,  $\rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S)$ ,  $\rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S)$  are equal in all but above two components  $i_2$  and  $i_2$ . Thus,

$$\begin{aligned} \left( \rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S) \right)^{\leq i} &= \left( \rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S) \right)^{\leq i} \cup \{s_1, s_2\} \\ \left( \rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S) \right)^{\leq i} &= \left( \rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S) \right)^{\leq i} \cup \{s_1\} \end{aligned}$$



$$\left(\rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S)\right)^{\leq i} = \left(\rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S)\right)^{\leq i} \cup \{s_2\}$$

By Lemma E.2,

$$g_i\left(\left(\rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S)\right)^{\leq i}\right) - g_i\left(\left(\rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S)\right)^{\leq i}\right) - g_i\left(\left(\rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S)\right)^{\leq i}\right) + g_i\left(\left(\rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S)\right)^{\leq i}\right) \leq 0$$

Then  $\frac{\partial^2 f}{\partial x_{i_1, s_1} \partial x_{i_2, s_2}}(\mathbf{x}) \leq 0$ .

**Case 3 :**  $i_1 \neq i_2, S_{i_1} \notin \{s_1, \circ\}$  or  $S_{i_2} \notin \{s_2, \circ\}$ : If  $S_{i_1} \notin \{s_1, \circ\}$ , then  $\rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S) = \rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S)$  and  $\rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S) = \rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S)$ . If  $S_{i_2} \notin \{s_2, \circ\}$ , then  $\rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S) = \rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S)$  and  $\rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S) = \rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S)$ . Either way, we have,

$$g_i\left(\left(\rho_{\vee, \vee}^{i_1, s_1, i_2, s_2}(S)\right)^{\leq i}\right) - g_i\left(\left(\rho_{\vee, \wedge}^{i_1, s_1, i_2, s_2}(S)\right)^{\leq i}\right) - g_i\left(\left(\rho_{\wedge, \vee}^{i_1, s_1, i_2, s_2}(S)\right)^{\leq i}\right) + g_i\left(\left(\rho_{\wedge, \wedge}^{i_1, s_1, i_2, s_2}(S)\right)^{\leq i}\right) = 0.$$

Then  $\frac{\partial^2 f}{\partial x_{i_1, s_1} \partial x_{i_2, s_2}}(\mathbf{x}) = 0$ .

In all cases,  $\frac{\partial f}{\partial x_{i_1, s_1} \partial x_{i_2, s_2}}(\mathbf{x}) \leq 0$ , thus  $f(\mathbf{x})$  is DR-submodular.  $\square$

**Corollary 5.6.** *There is an algorithm for attaining the expected  $(1 - 1/e)$ -regret of*

$$\mathcal{R}_{1-1/e}(T) \leq O\left(\left(|G|\right)^{10/3} T^{2/3} \log T\right)$$

on any  $(S_{OL}, \mathcal{G}_{SS})$ -bandit.

*Proof of Corollary 5.6.* Since EXT<sub>SS</sub> satisfies the conditions in Lemma 5.1 and the dimension of  $\mathcal{K}$  is  $d = |G|(|G| - 1) = O(|G|^2)$ . This is a direct corollary of Lemma 5.1.  $\square$

## F. Remark on the Stochastic Submodular Bandit

In this section, we show how our algorithms for adversarial setting can be applied to the stochastic setting proposed by Nie et al. (2022). We take the stochastic monotone submodular bandit with cardinality constraint investigated in (Nie et al., 2022) as an example.

**Stochastic submodular bandit model** In the stochastic model, there is an unknown distribution  $\mathcal{D}$ , its support is a set of set functions, we denote the set  $\text{supp}(\mathcal{D})$ . Any set function  $g \in \text{supp}(\mathcal{D})$  is defined on the power set of the ground set  $G$ , mapping a subset of  $G$  to a reward between  $[0, 1]$ , that is,  $g : 2^G \rightarrow [0, 1]$ . In  $t$ -th round, the reward function  $g'_t$  is drawn from  $\mathcal{D}$  and we can only select a subset  $S_t \subseteq G$  such that  $|S_t| \leq k$ , which is a cardinality constraint. Then we gain reward  $g'_t(S_t)$ . Note that the model does not need  $g'_t$  to be a monotone submodular function, but requires  $g = \mathbb{E}_{g'_t \sim \mathcal{D}}[g'_t]$  to be monotone submodular. Our goal is to minimize the  $(1 - 1/e)$ -regret:

$$\mathcal{R}_{1-1/e}^{sto}(T) = \left(1 - \frac{1}{e}\right)T \cdot \max_{S^* \subseteq G, |S^*| \leq k} g(S^*) - \mathbb{E} \left[ \sum_{t=1}^T g'_t(S_t) \right] = \left(1 - \frac{1}{e}\right)T \cdot \max_{S^* \subseteq G, |S^*| \leq k} g(S^*) - \mathbb{E} \left[ \sum_{t=1}^T g(S_t) \right]$$

The last equality is because the randomness of  $S_t$  is independent of the randomness of  $g'_t$ .

**Apply our algorithm on stochastic bandit model** Since cardinality constraint is a special case of the partition matroid constraint, we use our algorithm in Section 5.2. While applying our algorithm on the stochastic model, we see  $g_t = g = \mathbb{E}_{g'_t \sim \mathcal{D}}[g'_t]$  as the online function selected by the adversary, thus the online functions are monotone submodular. Note that, to obtain a regret bound w.r.t. the online reward function  $\{g_t\}_{t=1}^T$ , we need to query the value of  $g_t$  at some subset  $S_t$  in round  $t$ . However, since the algorithm is actually running on the stochastically realized function sequence  $\{g'_t\}_{t=1}^T$ , if we query the function value of  $S_t$ , the feedback is  $g'_t(S_t)$  rather than  $g_t(S_t) = g(S_t)$ . Fortunately, this is not a big issue since  $g'_t(S_t)$  is an unbiased estimate of  $g_t(S_t)$ . Now we go back to the proof of Lemma 5.1, and replace all the  $g_{t_q}(S_{t_q}) = g(S_{t_q})$  with  $g'_{t_q}(S_{t_q})$ . Since  $g'_{t_q}(S_{t_q}) \leq 1$ , it won't affect our bound for the dual local norm of the gradient

estimator. And since the randomness of the stochastic function  $g'_{t_q}$  is independent of all the randomness introduced in our algorithm and  $\mathbb{E}[g'_{t_q}(S_{t_q})] = g_{t_q}(S_{t_q})$ , the new gradient estimator constructed by replacing  $g_{t_q}(S_{t_q})$  with  $g'_{t_q}(S_{t_q})$  is still an unbiased estimator. Thus, as the same as the proof of our algorithm for adversarial submodular bandit with partition matroid constraint, we have the same regret bound,

$$\mathcal{R}_{1-1/e}^{adv}(T) = \max_{S^* \subseteq G, |S^*| \leq k} \mathbb{E} \left[ \left(1 - \frac{1}{e}\right) \sum_{t=1}^T g(S^*) - \sum_{t=1}^T g(S_t) \right] \leq O((k|G|)^{5/3} T^{2/3} \log T)$$

That is,

$$\left(1 - \frac{1}{e}\right) T \cdot \max_{S^* \subseteq G, |S^*| \leq k} g(S^*) - \mathbb{E} \left[ \sum_{t=1}^T g(S_t) \right] \leq O((k|G|)^{5/3} T^{2/3} \log T)$$