

# Online Learning via Offline Greedy Algorithms: Applications in Market Design and Optimization

Rad Niazadeh

Chicago Booth School of Business, Operations Management, rad.niazadeh@chicagobooth.edu

Negin Golrezaei

MIT Sloan School of Management, Operations Management, golrezae@mit.edu

Joshua Wang

Google Research Mountain View, joshuawang@google.com

Fransisca Susan

MIT Sloan School of Management, Operations Management, fsusan@mit.edu

Ashwinkumar Badanidiyuru

Google Research Mountain View, ashwinkumarbv@google.com

Motivated by online decision-making in time-varying combinatorial environments, we study the problem of transforming offline algorithms to their online counterparts. We focus on offline combinatorial problems that are amenable to a constant factor approximation using a greedy algorithm that is robust to local errors. For such problems, we provide a general framework that efficiently transforms offline robust greedy algorithms to online ones using Blackwell approachability. We show that the resulting online algorithms have  $O(\sqrt{T})$  (approximate) regret under the full information setting. We further introduce a bandit extension of Blackwell approachability that we call Bandit Blackwell approachability. We leverage this notion to transform greedy robust offline algorithms into a  $O(T^{2/3})$  (approximate) regret in the bandit setting. Demonstrating the flexibility of our framework, we apply our offline-to-online transformation to several problems at the intersection of revenue management, market design, and online optimization, including product ranking optimization in online platforms, reserve price optimization in auctions, and submodular maximization. We also extend our reduction to greedy-like first order methods used in continuous optimization, such as those used for maximizing continuous strong DR monotone submodular functions subject to convex constraints. We show that our transformation, when applied to these applications, leads to new regret bounds or improves the current known bounds. We complement our theoretical studies by conducting numerical simulations for two of our applications, in both of which we observe that the numerical performance of our transformations outperforms the theoretical guarantees in practical instances.

*Key words:* Blackwell approachability, Offline-to-online, No-regret, Submodular maximization, Product ranking, Reserve price optimization.

## 1. Introduction

We study the problem of designing efficient no-regret – also known as *vanishing regret* – online learning algorithms in complex real-world environments, where the underlying decision-making process is combinatorial in nature. In such environments, a decision-maker (learner) needs to experiment with exponentially many options whose rewards exhibit non-trivial and non-linear structures. Exploiting such structures to design efficient online learning algorithms is challenging as the underlying offline problems can indeed be NP-hard. Such offline problems can only admit approximation algorithms. Therefore, any efficient online learning algorithm can only hope to obtain vanishing regret with respect to an in-hindsight approximately optimal benchmark. This motivates our key research questions:

*How can one transform existing approximation algorithms for NP-hard offline problems to vanishing regret learning algorithms for a wide range of combinatorial environments? Can we efficiently exploit the combinatorial reward structure to eliminate the necessity of experimenting with exponentially many arms?*

To answer these questions, we consider an adversarial online learning setting. In every round  $t$ , the learner takes an action by choosing a (feasible) point  $\mathbf{z}_t$  among possibly exponentially many choices, and receives a reward of  $f_t(\mathbf{z}_t)$ . The adversarially chosen reward function  $f_t \in \mathcal{F}$ , which is unknown to the learner at the time of action, can be non-linear in action  $\mathbf{z}_t$ . We are interested in settings where the offline problem is NP-hard, and amenable to a  $\gamma$ -approximation algorithm, where  $\gamma \in (0, 1)$ .<sup>1</sup> In the offline problem, the reward function  $f \in \mathcal{F}$  is fully known, and the goal is to choose a feasible point  $\mathbf{z}$  that maximizes the obtained reward  $f(\mathbf{z})$ .

We focus on the prevalent class of offline approximation algorithms with a greedy nature. Roughly speaking, such approximation algorithms build up a solution stage by stage, choosing the next stage that offers the most local improvement with respect to a metric. We require the greedy approximation algorithms to be robust to local errors in every stage; for details, see Definition 5. Several combinatorial problems, ranging from classic submodular maximization problems to more recently studied optimization problems related to market design and revenue management, admit such robust greedy approximation algorithms. For details, see Section 6.

Our goal here is to conduct *offline-to-online transformations*; that is, to design online learning algorithms whose performance (over time) is as good as the performance of their corresponding offline approximation algorithm. The problem of offline-to-online transformation is studied by Kalai and Vempala (2005) and Dudík et al. (2017) when the learner can solve the offline problem efficiently. However, the approaches in these works fail when the learner only has access to an approximate

<sup>1</sup>Our framework can also be applied to polynomially solvable problems. In this case, the approximation factor is  $\gamma = 1$ .

solutions to the offline problem. This drawback is alleviated by Kakade et al. (2009) who study the offline-to-online transformation when (i) the offline problem is NP-hard but amenable to approximation, and (ii) the reward function is linear in the learner’s action. Kakade et al. (2009) crucially uses the linearity of the reward function (see also Garber (2021), Hazan et al. (2018)), and hence, their approach cannot be applied to our settings with nonlinear reward functions. We highlight that as shown by Hazan and Koren (2016), for a general offline problem, there may not exist an efficient offline-to-online transformation, justifying our assumption on the type of approximation algorithms.

We now summarize our main contributions.

**A framework for offline-to-online transformations.** We design a unified framework to transform robust greedy approximation algorithms to efficient online learning algorithms when the reward functions are not necessarily linear. We consider two online learning settings: *full information* and *bandit*. In the full information setting, the learner observes function  $f_t$  after taking action  $z_t$ , and in the bandit setting, the learner only observes the obtained reward  $f_t(z_t)$ .

For both settings, our proposed transformation relies on the celebrated Blackwell approachability theorem due to Blackwell (1956). The Blackwell approachability theorem is concerned with a two-player repeated game with a vector payoff, and presents a strategy under which the time-averaged vector payoff approaches some target set  $S$  that satisfies certain properties. As it is shown in Abernethy et al. (2011), there is a strong connection between Blackwell approachability and designing vanishing regret learning algorithms. In fact, for online linear optimization, they show that any strategy/algorithm for Blackwell approachability can be transformed to a vanishing regret learning algorithm and vice versa.

**Online learning algorithms using Blackwell strategies.** In this work, as one of our main contributions, we show that the transformation of Blackwell strategies to online vanishing regret algorithms is also possible for combinatorial non-linear learning settings whose underlying offline problem is NP-hard and admits a robust greedy  $\gamma$ -approximation algorithm. Specifically, we show that if the offline problem is *Blackwell reducible* (see Definitions 7 and 9), then we can design an online learning algorithm with vanishing  $\gamma$ -regret (as in Definition 2) by running a Blackwell algorithm for each stage (subproblem) of the offline greedy algorithm. In every round, these Blackwell algorithms are run sequentially to build up the learner’s action stage by stage. This allows the Blackwell algorithms to communicate with each other in a specific pattern dictated by the offline greedy algorithm. Thanks to such communication between Blackwell algorithms and the robustness of the offline greedy algorithm to local errors, the resulting online algorithm has a vanishing  $\gamma$ -regret. In fact, for the full information setting, we show that this transformation leads to an algorithm with  $O(N\sqrt{T})$   $\gamma$ -regret, where  $N$  is the number of subproblems in the offline algorithm.<sup>2</sup>

<sup>2</sup>Our regret bounds also depend on the diameter of vector payoff of the Blackwell games and their dimension; see Theorems 2 and 3. Further, in some applications we can show sub-linear dependency of the regret bound on the

The bandit setting turns out to be much trickier as the Blackwell algorithms cannot all obtain their desired feedback to update their course of actions over time. To circumvent this obstacle, we introduce a novel and customized bandit version of the Blackwell sequential game that we call *bandit Blackwell*. In this version, the player/algorithm does not obtain any feedback on his payoff unless he agrees to pay a certain cost. When the player agrees to pay such a cost, an extra “exploration” will be done, and he obtains an unbiased estimator of his payoff. Surprisingly, we show that in the bandit Blackwell sequential games, getting a vanishing regret with respect to a combination of approachability and exploration cost minimization is feasible (Theorem 3). We further give a tight lower bound on the rate of convergence for bandit Blackwell sequential games (Theorem 7).

Leveraging our notions of bandit Blackwell sequential games and approachability, we present an offline-to-online transformation in which  $N$  bandit Blackwell algorithms communicate with each other to build up a solution. To mimic the extra exploration step of bandit Blackwell games, we show how this communication can be interrupted in a controlled way when one of the bandit algorithms requests acquiring feedback. We also show how the required unbiased estimator of the vector payoff can be constructed. These pieces give us an online algorithm with  $O(NT^{2/3})$   $\gamma$ -regret.

**Applications.** Finally, to demonstrate the generality and effectiveness of our framework, we apply our offline-to-online transformation to several problems at the intersection of revenue management, market design, and online optimization that have been proposed and studied in the literature. In particular, we consider problems of (i) optimizing product ranking, (ii) optimizing personalized reserve prices in second price auctions, and (iii) submodular maximization (SM) in discrete and continuous domains (see Table 1). We show that in most cases, our transformations lead to new or improved regret bounds. In the following, we discuss our bounds in detail.

• **Product ranking optimization.** Online marketplaces have the opportunity of optimizing the ranking of displayed products in order to improve revenue, shape the demand, and reduce users’ search cost (see, for example, [Athey and Ellison \(2011\)](#), [Kempe et al. \(2003\)](#), [Ursu \(2016\)](#), [Aouad and Segev \(2021\)](#), [Derakhshan et al. \(2020\)](#), and [Golrezaei et al. \(2021b\)](#)). Inspired by this, we study the product ranking problem in the online adversarial setting. In this problem, the platform needs to identify a ranking/permutation of  $n$  items across  $n$  positions where items placed in top positions (positions with lower indices) get more visibility. The goal of the platform is maximize its user engagement (also known as market share), which is the probability that a consumer does not leave the platform without taking a desired action. To express user engagement as a function of the ranking over the products, we use the model proposed by [Asadpour et al. \(2020\)](#), which

number of sub problems, which turns out to be crucial for transforming continuous optimization algorithms to their online variants; see Theorems 9 and 10 in Appendix F.

is a generalization of the model presented by [Ferreira et al. \(2021\)](#). Under this model, the offline ranking problem can be written as maximizing sequential submodular functions; see [Section 6.1](#) for the definition of these functions. By applying our framework to this problem, we get  $O(n\sqrt{T\log n})$  and  $O(n^{5/3}(\log n)^{1/3}T^{2/3})$   $\frac{1}{2}$ -regret in full information and bandit settings, respectively. We note that our work is the first one that studies the product ranking problem under the aforementioned model in an online adversarial setting.<sup>3</sup>

• **Optimizing personalized reserve prices.** Second price auctions with reserve prices are prevalent in many marketplaces including online advertising markets, making them objects of both wide practical relevance and scientific interest (see, for example, [Hartline and Roughgarden \(2009\)](#), [Cesa-Bianchi et al. \(2014\)](#), [Beyhaghi et al. \(2018\)](#), [Roughgarden and Wang \(2019\)](#), [Golrezaei et al. \(2021a\)](#)). We study the online problem of optimizing personalized reserve prices, where buyers' valuations are chosen adversarially in every round. In the offline version of this problem, a seller wants to sell an item to one of  $n$  bidders by running a second price auction with personalized reserve prices. Each bidder  $i$  has a private value for the item. The seller wishes to maximize his revenue by optimizing over bidders' reserve prices. By applying our framework to the offline greedy algorithm of [Roughgarden and Wang \(2019\)](#), we achieve  $O(n\sqrt{T\log T})$   $\frac{1}{2}$ -regret in the full-information setting and  $O(n^{3/5}T^{4/5}(\log nT)^{1/3})$   $\frac{1}{2}$ -regret in the bandit setting. Our results match the previous bound for the full-information setting by [Roughgarden and Wang \(2019\)](#) who apply a slight variant of the Follow-the-Perturbed-Leader algorithm of [Kalai and Vempala \(2005\)](#) every round for each bidder; the bandit setting had not been studied prior to our work.<sup>4</sup>

• **Submodular maximization problems.** Many optimization problems that arise in the real world, including revenue management problems, can be expressed as maximizing a submodular function. The notion of submodularity is commonly used to describe the diminishing return property in discrete and continuous domains. Examples include the welfare maximization problem (e.g., [Dobzinski and Schapira \(2006\)](#) and [Vondrák \(2008\)](#)), capital budgeting with risk-averse investors (e.g., [Weingartner \(1967\)](#) and [Ahmed and Atamtürk \(2011\)](#)), and the problem of maximizing influence through the network (e.g., [Kempe et al. \(2003\)](#)).

We apply our framework to the adversarial online submodular maximization problem. For the online problem of maximizing monotone set submodular functions subject to cardinality constraints

<sup>3</sup> The offline PAC learning problem which resembles aspects of the online learning in the stochastic setting, is studied by [Ferreira et al. \(2021\)](#) for a special case of our model. PAC stands for probably approximately correct.

<sup>4</sup> In the special case with symmetric buyers and uniform reserve prices (also known as anonymous reserve auction, cf. [Alaei et al. \(2019\)](#)), minimizing regret under stochastic bandit setting is studied in [Cesa-Bianchi et al. \(2014\)](#), in which they obtain  $O(n\sqrt{T})$  regret bound. Here, the offline problem of finding the uniform optimal reserve can be solved exactly in polynomial time.

with size  $k$ , we transform the offline greedy algorithm by [Nemhauser et al. \(1978\)](#), which is a  $(1 - 1/e)$ -approximation, to yield  $O(k\sqrt{T\log n})$   $(1 - 1/e)$ -regret in the online full-information setting, matching the bound by [Streeter and Golovin \(2008\)](#) who use a variation of the EXP3 algorithm. Furthermore, our framework gives  $O(kn(\log n)^{1/3}T^{2/3})$   $(1 - 1/e)$ -regret in the bandit setting, improving the previous bound of  $O(k^2(n\log n)^{1/3}T^{2/3}(\log T)^2)$   $(1 - 1/e)$ -regret by [Streeter and Golovin \(2008, 2007\)](#) in the opaque feedback model, which is the limited feedback model that is analog to our bandit feedback model under exploration. See [Section 5](#) for more details.

For the online problem of maximizing non-monotone set submodular functions without any constraints, we transform a variation of the bi-greedy offline algorithm by [Buchbinder and Feldman \(2018\)](#) using our framework and obtain  $O(nT^{1/2})$   $\frac{1}{2}$ -regret in the full-information setting, matching the previous bound by [Roughgarden and Wang \(2018\)](#) who also take advantages of the bi-greedy offline algorithm of [Buchbinder and Feldman \(2018\)](#). Here,  $n$  is the number of coordinates. For the bandit setting, our transformation yields  $O(nT^{2/3})$   $\frac{1}{2}$ -regret. To the best of our knowledge, this is the first regret bound for the bandit setting of this challenging problem.

Switching to continuous submodular maximization settings, for the online problem of maximizing non-monotone continuous submodular functions without any constraints, we transform a variation of the continuous bi-greedy algorithm by [Niazadeh et al. \(2018\)](#) and obtain  $O(n\sqrt{T\log T})$   $\frac{1}{2}$ -regret in the online full-information setting. For the bandit setting, we obtain  $O(nT^{4/5}(\log T)^{1/3})$   $\frac{1}{2}$ -regret when the continuous submodular functions is weak-DR.<sup>5</sup> Our results for weak-DR submodular functions trivially yield results for strong-DR submodular functions. We highlight that the notion of weak-DR submodularity is equivalent to continuous submodularity and is easier to satisfy than strong-DR submodularity, which additionally requires coordinate-wise concavity; see the definition of weak-DR and strong-DR submodular functions in [Section E](#) in the appendix. Our work is the first one that designs online algorithms for weak-DR submodular functions. Furthermore, our bounds improve the previous bounds for strong-DR submodular functions by [Thang and Srivastav \(2021\)](#), which are  $O(T^{5/6})$   $\frac{1}{4}$ -regret and  $O(T^{11/12})$   $\frac{1}{4}$ -regret for the full-information and bandit settings, respectively.

For the problem of maximizing monotone continuous strong-DR submodular functions over a downward closed bounded convex set, by applying our framework to a variant of the Frank-Wolfe algorithm in [Bian et al. \(2017\)](#), for the full information setting, we design an online algorithm with  $O(\sqrt{Tn\log n})$   $(1 - 1/e)$ -regret. In terms of dependency on  $T$ , our regret bound matches the best regret bound in the literature by [Chen et al. \(2018b\)](#), which is also obtained by an online learning

<sup>5</sup> We omit the dependence on the Lipschitz constant here.

algorithm based on the Frank-Wolfe idea.<sup>6</sup> For the bandit setting, we design an algorithm with  $O(n(\log n)^{1/6}T^{5/6})$   $(1 - 1/e)$ -regret, improving the previous bound in Zhang et al. (2019), which is  $O(nT^{8/9})$  for the same approximation factor. See Theorems 9 and 10 in Appendix F.

**Experiments.** To demonstrate the practicality and ease-of-use of our framework, we evaluate our online learning algorithms for the product ranking and maximizing multiple reserves applications numerically. For both applications, our frameworks do better than the benchmark for both full-information and bandit settings on average. Furthermore, as expected, the full-information algorithm has smaller cumulative regret compared to the bandit algorithm. More details on the experiment is in Section A in the appendix.

**Table 1** Our results for selective applications of our framework, compared to previously known results.

Application	Approx Factor ( $\gamma$ )	Online Full-Information Setting		Online Bandit Setting	
		Our $\gamma$ -Regret Bound	The Best Prior Bound	Our $\gamma$ -Regret Bound	The Best Prior Bound
Product Ranking Problem	1/2	$O(n\sqrt{T}\log n)$	-	$O(n^{5/3}(\log n)^{1/3}T^{2/3})$	-
Reserve Price Optimization	1/2	$O(n\sqrt{T}\log T)$	$O(n\sqrt{T}\log T)$ *	$O(n^{3/5}T^{4/5}(\log nT)^{1/3})$	-
Monotone Set SM with Cardinality Constraints	$1 - 1/e$	$O(k\sqrt{T}\log n)$	$O(k\sqrt{T}\log n)$ †	$O(kn^{2/3}(\log n)^{1/3}T^{2/3})$	$O(k^2(n\log n)^{1/3}T^{2/3}(\log T)^2)$ †
Non-Monotone Set SM Functions	1/2	$O(n\sqrt{T})$	$O(n\sqrt{T})$ ‡	$O(nT^{2/3})$	-
Non-monotone Continuous SM (Strong-DR) Functions	1/2	$O(n\sqrt{T}\log T)$	$\gamma = 1/4, O(T^{5/6})$ §	$O(nT^{4/5}(\log T)^{1/3})$	$\gamma = 1/4, O(T^{11/12})$ §
Non-monotone Continuous SM (Weak-DR) Functions	1/2	$O(n\sqrt{T}\log T)$	-	$O(nT^{4/5}(\log T)^{1/3})$	-
Monotone Cont. SM (Strong-DR) in Downward Closed Convex Set	$1 - 1/e$	$O(\sqrt{Tn}\log n)$	$O(\sqrt{T})$ ¶	$O(n(\log n)^{1/6}T^{5/6})$	$O(nT^{8/9})$ ¶

\* Roughgarden and Wang (2019) † Streeter and Golovin (2008); ‡ Roughgarden and Wang (2018); § Thang and Srivastav (2021); ¶ Chen et al. (2018b);  
 ¶ Zhang et al. (2020);

## 1.1. Further Related Work

*Combinatorial learning.* Our work is related to the literature on online combinatorial learning. While in our work we study the design of efficient online learning algorithms for combinatorial problems whose loss function is not necessarily linear in the chosen action, the work on combinatorial learning

<sup>6</sup> Chen et al. (2018b), however, considers maximizing monotone continuous strong-DR submodular functions over a convex set which may not be downward closed. See also Thang and Srivastav (2021) for a work that builds on the Frank-Wolfe algorithm in Chen et al. (2018b) for the non-monotone strong-DR submodular maximization in a downward closed convex set. They obtain  $O(T^{3/4})$   $1/e$ -regret for the full-information setting.

mostly focuses on linear loss functions; see, for example, Abernethy et al. (2008), Uchiya et al. (2010), Cesa-Bianchi and Lugosi (2012), Audibert et al. (2014), Chen et al. (2013), Combes et al. (2015), and Zimmert et al. (2019). This line of work examines both the full-information and bandit settings. The standard exponentially weighted average forecaster obtains a tight  $O\left(m\sqrt{T\log\frac{d}{m}}\right)$  regret in the full-information setting, where  $m$  is the maximum  $\ell_1$ -norm of action vectors (Audibert et al. (2014)). The state-of-the-art regret bound for the bandit setting is  $O\left(\sqrt{dm^3T\log\frac{d}{m}}\right)$ , as reported in several papers (Bubeck et al. (2012), Cesa-Bianchi and Lugosi (2012), Hazan and Karnin (2016)). Our framework achieves matching regret with respect to  $T$  in the full-information setting without requiring the loss function to be linear. We get a worse regret (proportional to  $T^{2/3}$ ) for the bandit setting to account for the non-linearity in loss functions.

*Online adversarial submodular optimization.* In the previous section, we briefly mentioned some of the work that are closely related to our results on maximizing submodular functions in an online adversarial setting. Here, we provide more details. Chen et al. (2018a, 2020) use method based on Frank-Wolfe to design vanishing-regret learning algorithms for maximizing monotone continuous strong-DR submodular functions with convex constraints. Chen et al. (2018a) (respectively Chen et al. (2020)) assume that the algorithm can access to  $T^{1/2}$  exact (respectively  $T^{2/3}$  stochastic) gradient evaluations in every round and design an algorithm whose  $(1 - 1/e)$ -regret is  $O(\sqrt{T})$ .<sup>7</sup> The results of Chen et al. (2018a, 2020) were later improved by Zhang et al. (2019) who design another Frank-Wolfe inspired learning algorithm that has access to one stochastic gradient in each round and obtains  $O(T^{4/5})$   $(1 - 1/e)$ -regret. Zhang et al. (2019) further present a learning algorithm in the bandit setting for the problem of maximizing monotone continuous strong-DR submodular functions subject to matroid constraints. Their algorithm obtain  $O(T^{8/9})$   $(1 - 1/e)$ -regret. To see how our framework partially improve these results, refer Table 1 for a detailed comparison with our results related to monotone continuous submodular maximization subject to downward closed convex sets, and also non-monotone continuous submodular maximization with box constraints (also known as unconstrained).

*Online stochastic submodular optimization.* Designing learning algorithms for maximizing stochastic monotone continuous strong-DR submodular functions has been studied in Hassani et al. (2017), Mokhtari et al. (2020), Hassani et al. (2020), and Zhang et al. (2020). The best result for this setting is by Zhang et al. (2020) who obtain  $O(\sqrt{T})$   $(1 - 1/e)$ -regret using a stochastic variant of the Frank-Wolfe method. Their algorithm also implies the same regret bound for monotone set submodular maximization, which matches our regret bound for maximizing monotone set submodular function in the adversarial setting.

<sup>7</sup> The dependency on the number of elements  $n$  is not well specified in this work.

*Blackwell approachability.* Several aspects of Blackwell sequential game, including the design of efficient algorithms for Blackwell game with various information feedback structures, and the alternative conditions for approachability, have been studied in the literature. In terms of feedback structures, the original Blackwell game develops efficient projection algorithm for games that return the adversary’s moves on each round. Mannor et al. (2011) develop simple and efficient algorithms for a variant of Blackwell game where on each round, the player only obtains a random signal whose distribution depends on the action of the player and the adversary (as opposed to the action of the adversary). This variant is called Blackwell approachability with partial monitoring, and is further studied in Mannor et al. (2014) and Kwon and Perchet (2017). In terms of equivalent conditions for approachability, aside from the original halfspace-satisfiability condition for approachability in Blackwell (1956), alternative conditions for approachability, including the response-satisfiability criteria that we use in this paper, can be found in Lehrer (2003), Vieille (1992), Spinat (2002), and Milman (2006).

Blackwell approachability has also been proven to be a quintessential tool in various applications, as shown in Even-Dar et al. (2009) and Mannor and Shimkin (2006). However, most applications do not involve NP-hard combinatorial problems, and use the best-fixed action in hindsight (no approximation factor) as the benchmark for regret. Furthermore, they only create one Blackwell instance on each round. In contrast, we create multiple Blackwell instances on each round because the problems we consider have combinatorial nature and can only be solved efficiently in multiple stages. Furthermore, since we are solving NP-hard combinatorial problems with an intractable offline problem, we use a  $\gamma$ -approximation benchmark in our regret.

**Organization.** In Section 2, we present the offline optimization problem, adversarial online learning framework, and Blackwell sequential games. Section 3 presents the offline greedy approximation algorithm. In Sections 4 and 5, we present our offline-to-online transformation in the full information and bandit settings, respectively. Section 6 provides our regret bounds for the product ranking problem and optimizing reserve prices. Our regret bounds for maximizing unconstrained non-monotone submodular functions and maximizing monotone submodular functions over a downward closed bounded convex set are respectively presented in Sections E and F in the appendix. In Section A in the appendix, we present our numerical studies.

## 2. Preliminaries and Notations

In this section, we formulate our adversarial online learning framework for approximation algorithms. We then give an overview of Blackwell approachability (Blackwell 1956), an important technical tool that we use in this paper.

## 2.1. Offline Optimization and Approximations

Let  $\mathcal{F}$  be a space of functions defined over a (discrete or continuous) domain  $\mathcal{D}$ . Assume that  $\mathcal{F}$  is closed under addition, i.e., for any two functions  $f_1, f_2 \in \mathcal{F}$ , we have  $f_1 + f_2 \in \mathcal{F}$ . In the *offline optimization problem*, the problem of interest is finding a point  $\mathbf{z}^* \in \mathcal{D}$  such that

$$\mathbf{z}^* \in \arg \max_{\mathbf{z} \in \mathcal{C}} f(\mathbf{z}), \quad (1)$$

where  $f : \mathcal{D} \rightarrow [0, 1]$ , which belongs to  $\mathcal{F}$ , is the objective function, and  $\mathcal{C} \subseteq \mathcal{D}$  is the feasible region.<sup>8</sup> We further denote the optimal objective value of problem (1) by OPT; that is,  $\text{OPT} = \max_{\mathbf{z} \in \mathcal{C}} f(\mathbf{z})$ . We focus on maximization problems in this paper, but our techniques and results can easily be extended to minimization problems as well.

We consider offline problems that are NP-hard to solve exactly, and at the same time are amenable to a  $\gamma$ -approximation algorithm for some constant  $\gamma \in (0, 1)$ .

**DEFINITION 1** ( $\gamma$ -APPROXIMATION OFFLINE ALGORITHM). An offline algorithm  $\mathcal{A}$  for problem (1) is a polynomial time  $\gamma$ -approximation algorithm if for every  $f \in \mathcal{F}$  returns a feasible (possibly randomized) point  $\hat{\mathbf{z}} \in \mathcal{C}$  in polynomial time in the size of the algorithm's input such that

$$\mathbb{E}[f(\hat{\mathbf{z}})] \geq \gamma \cdot \text{OPT}.$$

Here, the expectation is with respect to the randomness in algorithm  $\mathcal{A}$ . The constant  $\gamma \in (0, 1)$  is referred to as the *approximation factor* of algorithm  $\mathcal{A}$ .

## 2.2. Adversarial Online Learning and Approximations

*Framework.* In the adversarial online learning version of problem (1), there is a learner, denoted by ALG, who plays  $T$  rounds of a sequential game against an adversary, denoted by ADV. In each round  $t \in [T]$ , ADV picks a function  $f_t \in \mathcal{F}$  and simultaneously ALG takes an action by picking a feasible point  $\mathbf{z}_t \in \mathcal{C}$ . Then, ALG obtains a reward equal to  $f_t(\mathbf{z}_t)$  and receives a *feedback* concerning this round. We highlight that unlike the offline optimization Problem (1), the unknown function  $f_t$  is not observable to ALG when it chooses action  $\mathbf{z}_t$ , and he only knows that  $f_t$  belongs to  $\mathcal{F}$ . Furthermore, ALG picks his action at time  $t$  only given the feedback of previous rounds  $1, 2, \dots, t-1$ , and in that sense, ALG is an online learner. ALG's goal is to pick points  $\{\mathbf{z}_t\}_{t=1}^T$  given the feedback of each round to maximize the accumulated reward  $\sum_{t=1}^T f_t(\mathbf{z}_t)$  against a worst-case adversary ADV. In this paper, for the sake of brevity and simplicity, we limit our focus to worst-case oblivious adversaries, i.e., adversaries that pick the sequence  $f_1, f_2, \dots, f_T$  upfront.

<sup>8</sup> For maximization problems, which are the focus of this paper, we only need our functions to be upper bounded by a constant. However, for simplicity, we assume that our functions are upper bounded by one.

*Feedback structures.* We consider two feedback structures: (i) *full information feedback*, where ALG observes the entire function  $f_t$  after choosing  $\mathbf{z}_t$ , and (ii) *bandit feedback*, where ALG only observes the quantity  $f_t(\mathbf{z}_t)$  after choosing  $\mathbf{z}_t$ . Let  $\phi_t$  be the feedback that ALG receives after picking  $\mathbf{z}_t$ . Then, ALG's next action  $\mathbf{z}_{t+1}$  is a function of the history  $\mathcal{H}_t$ , where  $\mathcal{H}_t \triangleq \{(\mathbf{z}_1, \phi_1), \dots, (\mathbf{z}_t, \phi_t)\}$ . More formally, any learning algorithm ALG can be described as mappings  $\{\pi_{\text{ALG}}^{(t)}\}_{t=1}^T$ , where each  $\pi_{\text{ALG}}^{(t)}$  maps the history  $\mathcal{H}_{t-1}$  to action  $\mathbf{z}_t$  for any  $t \in [T]$ . The mapping  $\pi_{\text{ALG}}^{(t)}$  can be either deterministic or randomized.

*Benchmarks and regret.* We would like to design polynomial-time online learning algorithms for offline problems that are NP-hard to solve exactly. Thus, we use the adapted notion of *approximate regret* to quantify the performance of an online algorithm. This notion is the regret with respect to  $\gamma$  fraction of the objective value at the best in-hindsight point. The notion of  $\gamma$ -regret, which is formally defined below, is common in the literature, see, for example, [Kakade et al. 2009](#), [Dudík et al. 2017](#), and [Roughgarden and Wang 2018](#). At a high level, our goal is to take a  $\gamma$ -approximation offline algorithm, and transform it to an online algorithm ALG with a sublinear  $\gamma$ -regret.

DEFINITION 2 ( $\gamma$ -REGRET). Let  $\sigma = \{(\mathbf{z}_t, f_t)\}_{t=1}^T$  be a sequence of strategies realized by online learner ALG and adversary ADV. Then, for any such  $\sigma$  and  $\gamma \in (0, 1)$ ,  $\gamma$ -regret( $\sigma$ ) is defined as

$$\gamma\text{-regret}(\sigma) \triangleq \gamma \cdot \max_{\mathbf{z} \in \mathcal{C}} \sum_{t=1}^T f_t(\mathbf{z}) - \sum_{t=1}^T f_t(\mathbf{z}_t).$$

With a slight abuse of the notation, we denote the worst-case expected approximate regret of ALG against any (oblivious) adversary ADV as follows:

$$\gamma\text{-regret}(\text{ALG}) \triangleq \max_{\{f_t\}_{t=1}^T} \left\{ \mathbb{E}[\gamma\text{-regret}(\sigma)] : \sigma = \{(\mathbf{z}_t, f_t)\}_{t=1}^T, \mathbf{z}_t \in \mathcal{C} = \text{ALG's strategy at time } t \in [T] \right\},$$

where the expectation is with respect to any randomness in ALG.

### 2.3. Blackwell Sequential Games and Approachability

To transform offline approximation algorithms to efficient online learning algorithms, we take advantage of *Blackwell sequential games*. A Blackwell sequential game is a repeated two-player game characterized by a tuple  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$ . In this repeated game,  $\mathcal{X}$  and  $\mathcal{Y}$  are both compact convex sets representing the players' action spaces, and  $\mathbf{p} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$  is a biaffine vector payoff function.<sup>9</sup> Moreover, parameter  $d \in \mathbb{N}$  is known as the dimension of the Blackwell sequential game. The vector payoff function  $\mathbf{p}$  is assumed to be known by both players. The game is played in  $T$  rounds. Each round involves player 1 choosing an action  $\mathbf{x}_t \in \mathcal{X}$  and player 2 choosing an action  $\mathbf{y}_t \in \mathcal{Y}$  simultaneously. Both actions may depend on the observed history  $((\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_{t-1}, \mathbf{y}_{t-1}))$ . This

<sup>9</sup> Function  $\mathbf{p}(\cdot, \cdot)$  is biaffine if for any  $\mathbf{x} \in \mathcal{X}$ ,  $\mathbf{p}(\mathbf{x}, \cdot)$  is affine and for any  $\mathbf{y} \in \mathcal{Y}$ ,  $\mathbf{p}(\cdot, \mathbf{y})$  is affine.

pair of actions produces the vector payoff  $\mathbf{p}(\mathbf{x}_t, \mathbf{y}_t)$ . The objective of player 1 is to ensure that the time-averaged payoff approaches a closed and convex target set  $S \subseteq \mathbb{R}^d$ , and the objective of player 2 is to prevent this from happening.

DEFINITION 3 (BLACKWELL APPROACHABILITY). In the Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$ , a target set  $S$  is  $g(T)$ -approachable if there exists a player 1 strategy such that for every player 2's strategy, the resulting sequence of actions satisfies

$$d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) \leq g(T),$$

where for any vector  $\mathbf{w} \in \mathbb{R}^d$  and set  $S \subseteq \mathbb{R}^d$ ,  $d_\infty(\mathbf{w}, S) \triangleq \inf_{\mathbf{v} \in S} \|\mathbf{w} - \mathbf{v}\|_\infty$  is the  $\ell_\infty$ -distance of vector  $\mathbf{w}$  from set  $S$ .

In this paper, we focus on the  $\ell_\infty$  norm rather than the usual  $\ell_2$  norm since it is more suitable for our applications. Our bounds on the approachability term  $g(T)$  will depend on the scale of the problem, and more formally on the diameter  $D(\mathbf{p})$  of the payoff function  $\mathbf{p}$ , defined as

$$D(\mathbf{p}) \triangleq \max_{\mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}} \|\mathbf{p}(\mathbf{x}, \mathbf{y})\|_\infty. \quad (2)$$

Ideally, player 1 aims to develop a strategy so that the term  $g(T)$  in Definition 3 converges to 0 as  $T$  converges to  $+\infty$ , and hence would be able to approach the target set  $S$  asymptotically. However, not every closed and convex target set  $S$  is approachable. To help with characterizing which sets are approachable, we additionally define the concept of *response-satisfiability*.

DEFINITION 4 (RESPONSE-SATISFIABLE). A closed and convex target set  $S$  is response-satisfiable in the Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  if for every player 2's action  $\mathbf{y} \in \mathcal{Y}$ , there exists a player 1's action  $\mathbf{x} \in \mathcal{X}$  such that the vector payoff falls into the target set, that is  $\mathbf{p}(\mathbf{x}, \mathbf{y}) \in S$ .

Blackwell's landmark result (Blackwell 1956) is an equivalence of (asymptotic) approachability and response-satisfiability.<sup>10</sup> We extend this result in the following theorem.

THEOREM 1. A closed and convex target set  $S$  is  $O(D(\mathbf{p}) \log(d)^{1/2} T^{-1/2})$ -approachable in the Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  if and only if the set  $S$  is response-satisfiable, where  $D(\mathbf{p})$ , defined in Equation (2), is the  $\ell_\infty$  diameter of the payoff function  $\mathbf{p}$ , and  $d$  is the dimension of the game.

We present a detailed proof of Theorem 1 in Section B.2 in the appendix, which is an adaptation of the original result of Blackwell (1956). The main difference between the Blackwell's original result and Theorem 1 is how the distance between the average payoff and set  $S$  is computed. While

<sup>10</sup> There are other equivalent structural criteria for approachability similar to response-satisfiability; see Section B.1 in the appendix for a list of these conditions.

Blackwell uses norm 2, we apply norm infinity. To account for this difference, we use the equivalence between Blackwell approachability and online linear optimization (Abernethy et al. 2011). This equivalence allows us to apply regret bounds for the latter problem that uses an arbitrary norm to find new bounds for the approachability problem. The regret bounds (on online linear optimization) can then be obtained via using Follow-the-Regularized-Leader (Shalev-Shwartz et al. 2012) or Online Mirror Descent (Bubeck et al. 2015) algorithms.

We finish this subsection by a few remarks regarding our treatment of the Blackwell approachability.

REMARK 1. As our goal is designing polynomial-time online learning algorithms, we further use algorithmic results in Even-Dar et al. 2009, and Abernethy et al. (2011) due to the equivalence between Blackwell approachability and full information adversarial online linear optimization. These results provide a polynomial-time approachable online algorithm satisfying the bound in Theorem 1, given access to a separation oracle for the closed and convex set  $S$ .<sup>11</sup> From this point on, when set  $S$  is response-satisfiable, we assume access to such an online algorithm that uses a separation oracle for the convex set  $S$  in a blackbox fashion.

REMARK 2. Another upshot of the above line of research on the equivalence between Blackwell approachability and full information online linear optimization is that an algorithm for player 1 to approach set  $S$  might only have access to the realized vector payoffs  $(\mathbf{p}(\mathbf{x}_1, \mathbf{y}_1), \dots, \mathbf{p}(\mathbf{x}_{t-1}, \mathbf{y}_{t-1}))$  in round  $t$ , rather than the entire history  $((\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_{t-1}, \mathbf{y}_{t-1}))$ , and this is indeed without loss of generality for obtaining the optimal bound of Theorem 1 (Abernethy et al. 2011). We relax this assumption in our “bandit Blackwell sequential game”, where we assume player 1 can only sometimes have access to an unbiased estimator of the realized vector payoff; see Section 5.1 for the definition and more details.

### 3. Approximation Algorithms for the Offline Problem: Iterative Greedy

As stated earlier, we are interested in transforming a  $\gamma$ -approximation algorithm for the offline problem (1) to an online learning algorithm, so that the worst-case  $\gamma$ -regret is sublinear in the number of rounds  $T$ . We consider a general class of algorithms for obtaining such an approximation guarantee, named *Iterative Greedy (IG)* algorithms. In an algorithm in this class, roughly speaking, a sequence of locally optimal decisions with respect to a specific metric (which we elaborate on

<sup>11</sup> Given the separation oracle for convex set  $S$ , the running-time should be polynomial in  $d$ ,  $T$ , and the number of bits required to encode  $\mathcal{X}$ . We are also considering a computational model where either the realized vector payoff is given as feedback at the end of each round, or the vector payoff function  $\mathbf{p}$  can be evaluated efficiently at any given pair of actions  $(\mathbf{x}, \mathbf{y})$ .

more later) leads to picking the final point. This point then provably provides an approximation guarantee with respect to the global optimal solution of problem (1).

Formally, consider the following abstract skeleton. Suppose that we have  $N$  subproblems indexed by  $i \in [N]$ . The algorithm starts from an initial feasible point  $\mathbf{z}^{(0)} \in \mathcal{C}$ . It then goes over the subproblems in the increasing order of their indices. The goal of each subproblem  $i$  is to return a new feasible point  $\mathbf{z}^{(i)} \in \mathcal{C}$  given the output of the previous subproblem, i.e.,  $\mathbf{z}^{(i-1)}$ . The algorithm finishes by returning the point  $\mathbf{z}^{(N)}$ . Now, each subproblem  $i$  performs two steps:

1. Local optimization: We associate a space of *update parameters*  $\Theta \subseteq \mathbb{R}^{d_{\text{param}}}$  to each subproblem. Given the previous point  $\mathbf{z}^{(i-1)}$  and the objective function  $f$ , the goal of this step is to find a *locally optimal* update parameter  $\boldsymbol{\theta}^{(i)} \in \Theta$  that satisfies:

$$\text{PAYOFF}(\boldsymbol{\theta}^{(i)}, \mathbf{z}^{(i-1)}, f) \geq \mathbf{0},$$

where  $\text{PAYOFF} : \Theta \times \mathcal{D} \times \mathcal{F} \rightarrow \mathbb{R}^{d_{\text{payoff}}}$  denotes the *parameter vector payoff function*.

2. Local update: Given the update parameter  $\boldsymbol{\theta}^{(i)}$  and  $\mathbf{z}^{(i-1)}$ , this step returns the next point

$$\mathbf{z}^{(i)} = \text{LOCAL-UPDATE}(\boldsymbol{\theta}^{(i)}, \mathbf{z}^{(i-1)}) \in \mathcal{C}.$$

Notably, we allow  $\text{LOCAL-UPDATE} : \Theta \times \mathcal{D} \rightarrow \Delta(\mathcal{C})$  to incorporate randomness, and therefore  $\mathbf{z}^{(i)}$  can be potentially a randomized point.

The above procedure is summarized in Algorithm 1.

REMARK 3. To simplify the notation, we only consider symmetric subproblems in this section, i.e., all of the subproblems have the same update parameter spaces, local optimization steps, etc. In some of our applications in Section 6, we need slightly different subproblems for different  $i = 1, \dots, N$ . Our method directly extends to that case by having index-dependent subproblems.

---

**Algorithm 1:** OFFLINE-IG( $\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta$ )

---

**Meta Input:** Feasible region  $\mathcal{C}$ , function space  $\mathcal{F}$ , defined over domain  $\mathcal{D}$ , parameter space

$\Theta \subseteq \mathbb{R}^{d_{\text{param}}}$ , and parameter vector payoff function  $\text{PAYOFF} : \Theta \times \mathcal{D} \times \mathcal{F} \rightarrow \mathbb{R}^{d_{\text{payoff}}}$ .

**Input:** function  $f \in \mathcal{F}$ .

**Output:** feasible point  $\mathbf{z} \in \mathcal{C}$ .

Initialize  $\mathbf{z}^{(0)} \in \mathcal{C}$ ; **for** subproblem  $i = 1$  to  $N$  **do**

    Choose update parameter  $\boldsymbol{\theta}^{(i)}$  so that  $\text{PAYOFF}(\boldsymbol{\theta}^{(i)}, \mathbf{z}^{(i-1)}, f) \geq \mathbf{0}$ ;

    Set  $\mathbf{z}^{(i)} \leftarrow \text{LOCAL-UPDATE}(\boldsymbol{\theta}^{(i)}, \mathbf{z}^{(i-1)})$ ;

Return the final point  $\mathbf{z} \leftarrow \mathbf{z}^{(N)}$ .

---

EXAMPLE 1. As a simple running example, consider the problem of maximizing a monotone submodular set function  $f : 2^{[n]} \rightarrow [0, 1]$  subject to the cardinality constraint  $k$ . A set function  $f : 2^{[n]} \rightarrow [0, 1]$  is submodular if for all  $S, T \subseteq [n]$ ,  $f(S \cup T) + f(S \cap T) \leq f(S) + f(T)$ . Maximizing a monotone submodular set function  $f : 2^{[n]} \rightarrow [0, 1]$  subject to the cardinality constraint  $k$  is an NP-hard optimization problem which admits the classic  $(1 - \frac{1}{e})$ -approximation greedy algorithm (Nemhauser et al. 1978). This algorithm starts from the empty set and picks elements greedily based on their marginal value to the current set, where the marginal value of adding element  $j$  to set  $S$  is  $f(S \cup \{j\}) - f(S)$ . That is, in each stage of this algorithm and given the chosen set so far  $S$ , an element  $j^* \in \arg \max_{j \in [n]} f(S \cup \{j\}) - f(S)$  with the highest marginal value is added to  $S$ . This problem is an example of problem (1) where  $\mathcal{D} = \{0, 1\}^n$ ,  $\mathcal{C} = \{\mathbf{z} \in \{0, 1\}^n : \mathbf{z} \cdot \mathbf{1}_n \leq k\}$  and  $\mathcal{F}$  is the space of all monotone submodular set functions. Here,  $\mathbf{1}_n$  is the all-ones vector with size  $n$ . The greedy algorithm is an instance of Algorithm 1 with  $\Theta = \Delta([n])$ , which is the set of all possible probability distributions over  $n$  elements, and  $N = k$  subproblems, one for each iteration of the greedy algorithm. To describe each subproblem, for  $\theta \in \Theta$ ,  $\mathbf{z} \in \mathcal{D}$ , and  $f \in \mathcal{F}$ ,

$$\forall j \in [n]: \quad [\text{PAYOFF}(\theta, \mathbf{z}, f)]_j = \theta^T \mathbf{y} - [\mathbf{y}]_j,$$

where  $[\cdot]_j$  denotes the  $j^{\text{th}}$  coordinate value of a vector and  $\mathbf{y} \triangleq [f(\mathbf{z} \cup \{j\}) - f(\mathbf{z})]_{j \in [n]}$  is the marginal objective value of adding element  $j$  to  $\mathbf{z}$ . Moreover,  $\text{LOCAL-UPDATE}(\theta, \mathbf{z})$  samples an element  $i^* \sim \theta$ , where  $\theta \in \Delta([n])$  is a probability distribution over  $n$  elements, and returns  $\mathbf{z} \cup \{i^*\}$ . Note that  $\text{PAYOFF}(\theta, \mathbf{z}, f) \geq \mathbf{0}$  guarantees  $\theta$  to only have positive mass on elements with maximum marginal value with respect to the point  $\mathbf{z}$ .

### 3.1. Approximation with Robustness to Local Errors

We focus on IG algorithms that (i) provide a worst-case multiplicative approximation guarantee for problem (1), and (ii) have a local optimization step that is robust to small errors, i.e., if we replace the locally optimal decisions with almost locally optimal ones, the final point still remains to be approximately optimal (with the same approximation factor), but up to a small additive error. The following definition formalizes this robustness notion.

DEFINITION 5 (( $\gamma, \delta$ )-ROBUST APPROXIMATION). An instance of Algorithm 1 is a  $(\gamma, \delta)$ -robust approximation algorithm for  $\gamma \in (0, 1)$  and  $\delta > 0$ , if it satisfies the following properties:

1. Algorithm 1 is a  $\gamma$ -approximation offline algorithm as in Definition 1,
2. Supposed that we replace  $\theta^{(i)}$  with  $\tilde{\theta}^{(i)}$  for every subproblem  $i = 1, \dots, N$ . Then, if

$$\forall j \in [d_{\text{payoff}}]: \quad \left[ \text{PAYOFF}(\tilde{\theta}^{(i)}, \mathbf{z}^{(i-1)}, f) \right]_j + \epsilon \geq 0,$$

then we should have:

$$\forall \hat{\mathbf{z}} \in \mathcal{C}: \quad \mathbb{E}[f(\hat{\mathbf{z}})] \geq \gamma \cdot f(\hat{\mathbf{z}}) - \delta N \epsilon,$$

where  $\epsilon > 0$  and  $[\cdot]_j$  denotes the  $j^{\text{th}}$  coordinate value of a vector.

For our purpose, we actually need a stronger version of this robustness property. This property essentially concerns multiple runs of the offline algorithm on a group of functions in  $\mathcal{F}$ , i.e.,  $\{f_t\}_{t \in [T]}$ , producing a sequence of feasible points  $\mathbf{z}_t \in \mathcal{C}$  for  $t \in [T]$ , and then guarantees a robust approximation for the summation function, i.e.,  $\sum_{t \in [T]} f_t(z)$ , against errors that are small on-average over these runs by the sequence  $\{\mathbf{z}_t\}_{t=1}^T$ . This property is satisfied in all of the applications that motivate our work, in particular in various set and continuous submodular maximization problems we study in Sections E and F in the appendix, and in both reserve price optimization and product ranking problems in Section 6.

**DEFINITION 6 (EXTENDED  $(\gamma, \delta)$ -ROBUST APPROXIMATION).** An instance of Algorithm 1 is an extended  $(\gamma, \delta)$ -robust approximation algorithm for  $\gamma \in (0, 1)$  and  $\delta > 0$ , if for any sequence of functions  $f_1, f_2, \dots, f_T \in \mathcal{F}$  the following property is satisfied:

- Consider a hypothetical variant of Algorithm 1 that in every round  $t$ , instead of choosing the update parameter  $\boldsymbol{\theta}_t^{(i)}$  for each subproblem  $i \in [N]$  such that  $\text{PAYOFF}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f) \geq \mathbf{0}$ , it chooses a “noisy” update parameter  $\tilde{\boldsymbol{\theta}}_t^{(i)}$  such that the sequence of update parameters  $(\tilde{\boldsymbol{\theta}}_t^{(i)})_{t \in [T]}$  satisfies the following condition

$$\forall j \in [d_{\text{payoff}}]: \left[ \sum_{t=1}^T \text{PAYOFF}(\tilde{\boldsymbol{\theta}}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) \right]_j + h(T) \geq 0.$$

Then, we have

$$\forall \hat{\mathbf{z}} \in \mathcal{C}: \sum_{t=1}^T \mathbb{E}[f_t(\mathbf{z}_t)] \geq \gamma \cdot \sum_{t=1}^T f_t(\hat{\mathbf{z}}) - \delta N h(T).$$

Here,  $\mathbf{z}_t^{(i)}$  is the output of subproblem  $i \in [N]$  for run  $t \in [T]$  by applying the hypothetical variant of Algorithm 1 on  $f_t$ ,  $h(\cdot): \mathbb{N} \rightarrow \mathbb{R}_+$ , and  $[\cdot]_j$  denotes the  $j^{\text{th}}$  coordinate value of a vector.

When there is only one run of the function (i.e.,  $T = 1$ ), the extended  $(\gamma, \delta)$ -robust approximation guarantee boils down to the weaker  $(\gamma, \delta)$ -robust approximation guarantee in Definition 5. We finish this section by revisiting our running example and demonstrating the (extended) robust approximation property in this example.

**EXAMPLE 1 (CONTINUED).** By digging deeper in the original analysis of the greedy algorithm (Nemhauser et al. 1978), we show that the greedy algorithm satisfies the extended  $(\gamma, \delta)$ -robust approximation property for  $\gamma = 1 - \frac{1}{e}$  and  $\delta = 1$ .

Suppose that  $\mathbf{z}^* = \{a_1, \dots, a_k\}$  is the optimal solution of the offline problem; that is,  $\mathbf{z}^* = \arg \max_{\mathbf{z} \in \{0,1\}^n: \mathbf{z}_1 \leq k} \sum_{t=1}^T f_t(\mathbf{z})$ . (We use binary indicator vectors and sets interchangeably in this paper.) Further, let  $\mathbf{z}_t^{(i)}$  be the solution returned by the  $i^{\text{th}}$  subproblem of the greedy algorithm when the objective function is  $f_t$ . Then, for every  $i \in [k]$ ,

$$\sum_{t=1}^T f_t(\mathbf{z}^*) - \sum_{t=1}^T f_t(\mathbf{z}_t^{(i-1)}) \stackrel{(1)}{\leq} \sum_{t=1}^T f_t(\mathbf{z}^* \cup \mathbf{z}_t^{(i-1)}) - \sum_{t=1}^T f_t(\mathbf{z}_t^{(i-1)})$$

$$\begin{aligned}
&= \sum_{t=1}^T \sum_{j=1}^k \left( f_t(\mathbf{z}_t^{(i-1)} \cup \{a_1, \dots, a_j\}) - f_t(\mathbf{z}_t^{(i-1)} \cup \{a_1, \dots, a_{j-1}\}) \right) \\
&\stackrel{(2)}{\leq} \sum_{j=1}^k \sum_{t=1}^T \left( f_t(\mathbf{z}_t^{(i-1)} \cup \{a_j\}) - f_t(\mathbf{z}_t^{(i-1)}) \right) \\
&\stackrel{(3)}{=} \sum_{j=1}^k \sum_{t=1}^T \left( \langle \tilde{\boldsymbol{\theta}}_t^{(i)}, \mathbf{y}_t^{(i-1)} \rangle - \left[ \text{PAYOFF}(\tilde{\boldsymbol{\theta}}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) \right]_{a_j} \right) \\
&= \sum_{j=1}^k \sum_{t=1}^T \left( \sum_{j=1}^n [\tilde{\boldsymbol{\theta}}_t^{(i)}]_j \left( f_t(\mathbf{z}_t^{(i-1)} \cup \{j\}) - f_t(\mathbf{z}_t^{(i-1)}) \right) - \left[ \text{PAYOFF}(\tilde{\boldsymbol{\theta}}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) \right]_{a_j} \right) \\
&= k \cdot \sum_{t=1}^T \mathbb{E} \left[ f_t(\mathbf{z}_t^{(i)}) - f_t(\mathbf{z}_t^{(i-1)}) \right] - \sum_{j=1}^k \sum_{t=1}^T \left[ \text{PAYOFF}(\tilde{\boldsymbol{\theta}}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) \right]_{a_j} \\
&\stackrel{(4)}{\leq} k \cdot \sum_{t=1}^T \mathbb{E} \left[ f_t(\mathbf{z}_t^{(i)}) - f_t(\mathbf{z}_t^{(i-1)}) \right] + kh(T),
\end{aligned}$$

where  $\mathbf{y}_t^{(i)} \triangleq \left[ f_t(\mathbf{z}_t^{(i)} \cup \{j\}) - f_t(\mathbf{z}_t^{(i)}) \right]_{j \in [n]}$ . In the above chain of inequalities, inequality (1) holds because function  $f_t$  is monotone, inequality (2) holds due to submodularity of functions  $\{f_t\}_{t=1}^T$ , equality (3) holds because of the definition of the payoff vector in Example 1, and inequality (4) holds because of the condition in Definition 6. By rearranging the terms and taking expectations, we have:

$$\sum_{t=1}^T \mathbb{E} \left[ f_t(\mathbf{z}^*) - f_t(\mathbf{z}_t^{(i)}) \right] \leq \left(1 - \frac{1}{k}\right) \sum_{t=1}^T \mathbb{E} \left[ f_t(\mathbf{z}^*) - f_t(\mathbf{z}_t^{(i-1)}) \right] + h(T).$$

By recursing the above inequality for  $i = 1, \dots, k$ , and rearranging the terms, we finally have:

$$\sum_{t=1}^T \mathbb{E} [f_t(\mathbf{z}_t)] \geq \left(1 - \left(1 - \frac{1}{k}\right)^k\right) \sum_{t=1}^T f_t(\mathbf{z}^*) - h(T) \sum_{i=1}^k \left(1 - \frac{1}{k}\right)^{i-1} \geq \left(1 - \frac{1}{e}\right) \sum_{t=1}^T f_t(\mathbf{z}^*) - kh(T).$$

■

Not all greedy algorithms have robust guarantees. Example 2 of Section C in the appendix shows why, e.g., Dijkstra's algorithm for the shortest path problem, is not robust to local errors.

#### 4. Online Algorithm under Full Information Feedback Structure

In this section, we show how to transform an offline IG algorithm (Algorithm 1) to an online learning algorithm with a small approximate regret whenever it (i) is an extended robust approximation algorithm (Definition 6), and (ii) satisfies an extra condition that we call *Blackwell reducibility*. We first introduce this condition. Then, with the help of the Blackwell approachability (Theorem 1), we propose a meta full information online learning algorithm as our offline-to-online transformation.

#### 4.1. Blackwell Reducibility

The crux of our technique to transform an offline IG algorithm to an online learning algorithm is the possibility of reducing the local optimization step of Algorithm 1 to an approachable instance of the Blackwell sequential game as in Section 2.3.

DEFINITION 7 (BLACKWELL REDUCIBILITY). An instance OFFLINE-IG( $\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta$ ) of Algorithm 1 is Blackwell reducible if there exists an instance  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  of the Blackwell sequential game (with a biaffine vector payoff function  $\mathbf{p}$ ) and a mapping AdvB :  $\mathcal{D} \times \mathcal{F} \rightarrow \mathcal{Y}$  called *synthetic Blackwell adversary function*, such that:

1. The player 1's action space  $\mathcal{X}$  is equal to the parameter space  $\Theta$  in Algorithm 1; i.e.,  $\mathcal{X} = \Theta$ , and for any  $\boldsymbol{\theta} \in \Theta, \mathbf{z} \in \mathcal{D}, f \in \mathcal{F}$ , we have  $\text{PAYOFF}(\boldsymbol{\theta}, \mathbf{z}, f) = \mathbf{p}(\boldsymbol{\theta}, \text{AdvB}(\mathbf{z}, f))$ .
2. The set  $S \triangleq \{\mathbf{u} \in \mathbb{R}^{d_{\text{payoff}}} : [\mathbf{u}]_j \geq 0, j \in [d_{\text{payoff}}]\}$  is response-satisfiable (Definition 4).

EXAMPLE 1 (CONTINUED). The greedy algorithm of Nemhauser et al. (1978) is Blackwell reducible. Consider an instance  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  of Blackwell where  $\mathcal{X} = \Theta = \Delta([n])$  and  $\mathcal{Y} = [0, 1]^n$ . The synthetic Blackwell adversary function is  $\text{AdvB}(\mathbf{z}, f) = [f(\mathbf{z} \cup \{j\}) - f(\mathbf{z})]_{j \in [n]}$ , and the biaffine Blackwell vector payoff function is  $\mathbf{p}(\boldsymbol{\theta}, \mathbf{y}) = \boldsymbol{\theta}^T \mathbf{y} \mathbf{1}_n - \mathbf{y}$ .<sup>12</sup> Recall that  $\mathbf{1}_n$  is all-ones  $n$ -dimensional vector. Furthermore, set  $S$  is response-satisfiable because for every player 2's action  $\mathbf{y} \in \mathcal{Y}$ , playing  $\boldsymbol{\theta} = e_{j^*}$  with  $j^* = \underset{j \in [n]}{\text{argmax}} y_j$  implies that  $\mathbf{p}(\boldsymbol{\theta}, \mathbf{y}) \geq \mathbf{0}$ .

#### 4.2. Offline-to-Online Transformation with Full Information Feedback

If the offline algorithm (Algorithm 1) is Blackwell reducible, then one can think of the following approach to transform it into an online learning algorithm: associate an instance of the Blackwell sequential game to each subproblem  $i$  following the Blackwell reducibility, and then running  $N$  parallel online approachable algorithms for these Blackwell instances to find a sequence of assignments of the update parameter of each subproblem  $i$  over time. We further need to show how to synchronize these parallel runs through a proper communication between them, so as to construct a sequence of feasible solutions  $\mathbf{z}_1, \dots, \mathbf{z}_T$  guaranteeing a small approximate regret.

4.2.1. **Overview of the Algorithm** Recall that our goal in the offline problem is to solve the optimization problem  $\max_{\mathbf{z} \in \mathcal{C}} f(\mathbf{z})$ , where  $f \in \mathcal{F}$ . The offline problem admits a polynomial time IG  $\gamma$ -approximation algorithm, OFFLINE-IG( $\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta$ ), presented in Algorithm 1. This algorithm solves  $N$  subproblems sequentially, building the solution step by step. In step/subproblem  $i$  of this algorithm, we first update parameters  $\boldsymbol{\theta}^{(i)} \in \Theta \subseteq \mathbb{R}^{d_{\text{param}}}$  using the previous point  $\mathbf{z}^{(i-1)}$ , and then return the next point  $\mathbf{z}^{(i)}$  to feed to the next subproblem. The algorithm finishes by returning the final point  $\mathbf{z}^{(N)}$ .

<sup>12</sup> Note that  $\mathcal{Y} = [0, 1]^n$  because  $f : 2^{[n]} \rightarrow [0, 1]$  is monotone non-decreasing.

As stated earlier, we assume that the offline problem is Blackwell reducible; that is, we can define the Blackwell instance  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  and synthetic Blackwell adversary function  $\text{AdvB} : \mathcal{D} \times \mathcal{F} \rightarrow \mathcal{Y}$  that satisfy the conditions in Definition 7. Although this definition might seem technical, verifying it for many offline algorithms is indeed straightforward; see Sections 6, E, and F.

For the online version of the above offline algorithm, the meta input is feasible region  $\mathcal{C}$ , function space  $\mathcal{F}$ , which is defined over domain  $\mathcal{D}$ , and parameter space  $\Theta \subseteq \mathbb{R}^{d_{\text{param}}}$ . We further consider having access to an online Blackwell algorithm  $\text{AlgB}$ , player 1's strategy in the above Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$ , where such algorithm (i) ensures that the distance between the average vector payoff  $\frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t)$  and set  $S$  goes to zero with rate  $g(T) = O\left(D(\mathbf{p}) \sqrt{\frac{\log(d)}{T}}\right)$  against any adversarial player 2's strategy (see Theorem 1), and (ii) can be implemented in polynomial time having access to a separation oracle for the convex set  $S$ . As stated earlier, the existence of such an algorithm follows from the work of Even-Dar et al. 2009 and Abernethy et al. (2011); see Remark 1. We consider  $N$  parallel copies of this algorithm, one for each subproblem  $i \in [N]$ . It is also important to note that in most of our applications, set  $S$  is the positive orthant, for which a polynomial time separation oracle exists.<sup>13</sup> Our algorithm that takes advantage of  $N$  parallel copies of the online Blackwell algorithm is summarized in Algorithm 2.

Let  $\text{AlgB}^{(i)}$  be the copy of the above online Blackwell algorithm associated to subproblem  $i \in [N]$ . This copy handles the local optimization step of subproblem  $i$  in the  $\text{OFFLINE-IG}(\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta)$  in every round  $t \in [T]$  without knowing function  $f_t$ . Consider the decision-making process of this online algorithm in round  $t$ . The inputs prior to this round are all the update parameters of the subproblem  $i$  in the first  $t - 1$  rounds, i.e.,  $\boldsymbol{\theta}_1^{(i)}, \dots, \boldsymbol{\theta}_{t-1}^{(i)}$ , and the realized vector payoffs of the first  $t - 1$  rounds against player 2 in the Blackwell sequential game associated to subproblem  $i$ , i.e.,  $\mathbf{p}(\boldsymbol{\theta}_1^{(i)}, \mathbf{y}_1^{(i)}), \dots, \mathbf{p}(\boldsymbol{\theta}_{t-1}^{(i)}, \mathbf{y}_{t-1}^{(i)})$ . We consider a particular player 2 for this Blackwell sequential game. More explicitly, the synthetic adversary function  $\text{AdvB}$ , which is part of our reduction, plays the role of player 2 in any round  $t$ , i.e.,  $\mathbf{y}_t^{(i)} = \text{AdvB}(z_t^{(i-1)}, f_t)$ . Given the input prior to time  $t$ ,  $\text{AlgB}^{(i)}$  returns the new update parameter  $\boldsymbol{\theta}_t^{(i)}$ .<sup>14</sup>

After the online Blackwell algorithm  $\text{AlgB}^{(i)}$  returns the update parameter  $\boldsymbol{\theta}_t^{(i)}$ , we return the point  $z_t^{(i)}$  by calling the  $\text{LOCAL-UPDATE}$  function in the offline algorithm, i.e., we set  $z_t^{(i)}$  to

<sup>13</sup> For the application of maximizing monotone Strong-DR submodular functions over downward closed bounded convex sets, presented in Section F, the target set  $S$  is not a positive orthant, yet it is convex and admits a polynomial separation oracle.

<sup>14</sup> Note that the adversary's action in round  $t$  for subproblem  $i$  is  $\text{AdvB}(z_t^{(i-1)}, f_t)$ , where  $z_t^{(i-1)}$  is the decision made by the first  $i - 1$  subproblems in round  $t$ . Here, both  $f_t$  and  $z_t^{(i-1)}$  can be viewed as the signals used by the adversary in round  $t$  to determine his strategy. Using these signals in the adversary's strategy (i.e.,  $\text{AdvB}(z_t^{(i-1)}, f_t)$ ) is allowed. This is because conditioned on the history of plays (and past signals) in previous rounds (i.e.,  $f_\tau$  and  $z_\tau^{(i-1)}$  for any  $\tau < t$ ), as well as the feedback that  $\text{AlgB}^{(i)}$  receives in this round (i.e.,  $\text{PAYOFF}(\boldsymbol{\theta}_t^{(i)}, z_t^{(i-1)}, f_t)$ ), both  $f_t$  and  $z_t^{(i-1)}$  will be independent of  $\text{AlgB}^{(i)}$ 's action in future rounds  $t' > t$ .

LOCAL-UPDATE( $\theta_t^{(i)}, z_t^{(i-1)}$ ). Observe that the point returned by the subproblem  $i$ , i.e.,  $z_t^{(i)}$ , depends on the point returned by the previous subproblem  $z_t^{(i-1)}$ . This highlights that while each online Blackwell algorithm is responsible for one subproblem, they communicate with each other to build the final solution, where this communication is structured by the offline algorithm through the LOCAL-UPDATE function. After obtaining the point  $z_t^{(i)}$ , we move to subproblem  $i + 1$ .

Finally note that simulating the actions of our particular player 2 to determine the realized vector payoffs of each round, and computing/sending this feedback at the end of each round to  $\text{AlgB}^{(i)}$  (as player 1) in a computationally efficient manner, require the following:

- Knowing the point  $z_t^{(i-1)}$  picked by subproblem  $i - 1$  at time  $t$ : This is possible as we go over our subproblems in the order  $i = 1, \dots, N$  in each round  $t$ .
- Knowing the function  $f_t$ : This is possible because here we study the full information feedback structure, where under this structure we have access to  $f_t$  after we choose point  $z_t = z_t^{(N)}$ .
- Being able to compute the realized vector payoff  $\mathbf{p}(\theta_t^{(i)}, \text{AdvB}(z_t^{(i-1)}, f_t))$  efficiently given  $\theta_t^{(i)}$ ,  $f_t$ , and  $z_t^{(i-1)}$ . This is possible as this quantity is equal to  $\text{PAYOFF}(\theta_t^{(i)}, z_t^{(i-1)}, f_t)$ , which can be evaluated in polynomial time as  $\text{OFFLINE-IG}(\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta)$  is a polynomial time algorithm.

**4.2.2. Regret Analysis** The following theorem, which bounds the regret of our algorithm, is the main result of this section.

**THEOREM 2 (Full information offline-to-online transformation).** *Suppose that an instance of the algorithm  $\text{OFFLINE-IG}(\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta)$  for the offline problem (1) satisfies the following properties:*

- *It is an extended  $(\gamma, \delta)$ -robust approximation for  $\gamma \in (0, 1)$  and  $\delta \in \mathbb{R}_+$ , as in Definition 6.*
- *It is Blackwell reducible, that is, we can define the Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  and synthetic Blackwell adversary function  $\text{AdvB}: \mathcal{D} \times \mathcal{F} \rightarrow \mathcal{Y}$  that satisfy the conditions in Definition 7.*

*Consider the full-information adversarial online learning version of the problem (1), and let  $\text{AlgB}$  be a polynomial time Blackwell algorithm for  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  as in Remark 1. Then, for this online problem,  $\text{ONLINE-IG}(\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta, \text{AlgB})$  runs in polynomial time and satisfies the following  $\gamma$ -regret bound:*

$$\gamma\text{-regret}(\text{ONLINE-IG}(\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta, \text{AlgB})) \leq O\left(D(\mathbf{p}) N \delta \sqrt{\log(d_{\text{payoff}} T)}\right),$$

*where  $N$  is the number of subproblems,  $d_{\text{payoff}}$  is the dimension of vector payoffs, and  $D(\mathbf{p})$ , defined in Equation (2), is the  $\ell_\infty$ -diameter of the vector payoff space.*

*Proof of Theorem 2.* Consider a subproblem  $i \in [N]$ . Let  $S$  be the  $d_{\text{payoff}}$ -dimensional positive orthant; see the Blackwell reducibility definition and its associated approachable set  $S$  in Definition 7. Because  $S$  is response-satisfiable and projection onto  $S$  can be done in polynomial-time,

**Algorithm 2:** Full-information Online Learning Meta-algorithm (ONLINE-IG)

**Meta Input:** Feasible region  $\mathcal{C}$ , function space  $\mathcal{F}$ , defined over domain  $\mathcal{D}$ , and parameter space  $\Theta \subseteq \mathbb{R}^{d_{\text{param}}}$ .

**Offline algorithm and reduction gadgets:** An instance OFFLINE-IG( $\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta$ ) of Algorithm 1, the Blackwell instance  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  and synthetic Blackwell adversary function AdvB:  $\mathcal{D} \times \mathcal{F} \rightarrow \mathcal{Y}$  as this offline algorithm is Blackwell reducible (Definition 7).

**Input:** Number of rounds  $T$ ; access to a Blackwell online algorithm AlgB.

**Output:** Points  $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_T \in \mathcal{C}$ .

Initialize  $N$  parallel instances  $\{\text{AlgB}^{(i)}\}_{i=1}^N$  of the online algorithm AlgB;

**for** round  $t = 1$  to  $T$  **do**

    Initialize  $\mathbf{z}_t^{(0)} \in \mathcal{C}$ ;

**for** subproblem  $i = 1$  to  $N$  **do**

        Choose update parameter  $\boldsymbol{\theta}_t^{(i)}$  by querying online algorithm AlgB<sup>(i)</sup> given the update parameters and vector payoffs prior to round  $t$  in the Blackwell sequential game of subproblem  $i$ , that is,  $\boldsymbol{\theta}_t^{(i)} \leftarrow \text{AlgB}^{(i)}\left(\boldsymbol{\theta}_1^{(i)}, \dots, \boldsymbol{\theta}_{t-1}^{(i)}, \mathbf{p}(\boldsymbol{\theta}_1^{(i)}, \mathbf{y}_1^{(i)}), \dots, \mathbf{p}(\boldsymbol{\theta}_{t-1}^{(i)}, \mathbf{y}_{t-1}^{(i)})\right)$ ;

        Set  $\mathbf{z}_t^{(i)} \leftarrow \text{LOCAL-UPDATE}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}) \in \mathcal{C}$ ;

**end**

    Play the final point  $\mathbf{z}_t \leftarrow \mathbf{z}_t^{(N)}$ ;

$\langle$  Full information feedback: adversary reveals function  $f_t \in \mathcal{F}$   $\rangle$ ;

**for**  $i = 1$  to  $N$  **do**

        Give feedback  $\mathbf{p}(\boldsymbol{\theta}_t^{(i)}, \mathbf{y}_t^{(i)}) \leftarrow \text{PAYOFF}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t)$  to the Blackwell Algorithm

        AlgB<sup>(i)</sup> (as the vector payoff of round  $t$  against player 2); // Note that

$\mathbf{y}_t^{(i)} = \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t)$  for player 2 implicitly, although we do not need to evaluate AdvB to compute this action explicitly.

**end**

**end**

there exists a polynomial-time online algorithm AlgB (with  $N$  parallel copies  $\{\text{AlgB}^{(i)}\}_{i=1}^N$ ) that guarantees Blackwell approachability for the Blackwell instance corresponding to subproblem  $i$  with  $g(T) = O\left(D(\mathbf{p}) \sqrt{\frac{\log(d_{\text{payoff}})}{T}}\right)$ , based on Theorem 1. Therefore, we have:

$$d_\infty\left(\frac{1}{T} \sum_{t=1}^T \mathbf{p}\left(\boldsymbol{\theta}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t)\right), S\right) = d_\infty\left(\frac{1}{T} \sum_{t=1}^T \mathbf{p}\left(\boldsymbol{\theta}_t^{(i)}, \mathbf{y}_t^{(i)}\right), S\right) \leq g(T).$$

Because the target set  $S$  is the positive orthant, we have

$$d_\infty\left(\frac{1}{T} \sum_{t=1}^T \mathbf{p}\left(\boldsymbol{\theta}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t)\right), S\right) \leq g(T) \iff \forall j : \left[ \sum_{t=1}^T \mathbf{p}\left(\boldsymbol{\theta}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t)\right) \right]_j \geq -Tg(T)$$

Because of Blackwell reducibility,  $\text{PAYOFF}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) = \mathbf{p}(\boldsymbol{\theta}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t))$ . Therefore,

$$\forall j \in [d_{\text{payoff}}] : \left[ \sum_{t=1}^T \text{PAYOFF}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) \right]_j \geq -Tg(T). \quad (3)$$

Finally, because Algorithm 1 is an extended  $(\gamma, \delta)$ -robust approximation (see Definition 6), from Equation (3), we have:

$$\sum_{t=1}^T \mathbb{E}[f_t(\mathbf{z}_t)] \geq \gamma \cdot \sum_{t=1}^T f_t(\mathbf{z}^*) - \delta NTg(T) = \gamma \cdot \sum_{t=1}^T f_t(\mathbf{z}^*) - O\left(\delta ND(\mathbf{p}) \sqrt{\log(d_{\text{payoff}})T}\right),$$

which finishes the proof. Here,  $\mathbf{z}^*$  is the optimal in-hindsight feasible solution, i.e.,  $\mathbf{z}^* = \operatorname{argmax}_{\mathbf{z} \in \mathcal{C}} \sum_{t=1}^T f_t(\mathbf{z})$ . ■

We finish this section by reviewing our running example (Example 1) and mentioning the regret bound we get as a direct corollary of Theorem 2.

EXAMPLE 1 (CONTINUED). The greedy algorithm in Nemhauser et al. (1978) is an extended  $(1 - \frac{1}{e}, 1)$ -robust approximation algorithm and Blackwell reducible. It has  $N = k$  subproblems, the  $\ell_\infty$  diameter of the payoff space is  $D = 1$ , and the dimension of vector payoffs is  $d = n$ . Therefore, by invoking Algorithm 2 given any Blackwell algorithm satisfying the approachability bound in Theorem 1, we obtain the following bound:

$$\left(1 - \frac{1}{e}\right)\text{-Regret}(\text{Algorithm 2}) \leq O(k\sqrt{\log(n)T}),$$

which exactly matches the bound known in Streeter and Golovin (2008) for the same problem.

## 5. Online Algorithm under Bandit Information Feedback Structure

So far, we presented a framework to transform an offline iterative greedy algorithm to its online counterpart under the full information feedback structure. While the full information setting provides the theoretical foundations for the rest of our results, from an application point of view, it is less motivated. In almost all applications of our framework in revenue management and online decision making (e.g., product ranking problem and reserve price optimization), assuming the learner has full information feedback is rather a strong assumption.

In this section, we seek to relax this assumption, and try to understand if our framework can be extended to the more challenging bandit feedback structure setting. Under the bandit feedback structure, at the end of each round  $t$ , the learner faces an additional challenge: he only has access to  $f_t(\mathbf{z}_t)$ , rather than the entire function  $f_t$  like in the full information setting. Such a feedback structure prevents the online Blackwell algorithms  $\text{AlgB}^{(i)}$  to receive the feedback they require.

To overcome this challenge, we first consider a stylized bandit variation of the sequential Blackwell game. We characterize a new notion of approachability that we call *bandit Blackwell approachability* and provide an algorithm achieving the information-theoretic tight approachability bound for this problem. This algorithm uses an algorithm for the full information version of the Blackwell sequential game in a blackbox fashion.

We then introduce the extra ingredient that is needed for our bandit transformation, which is the possibility of creating an unbiased estimator for the vector payoff of the Blackwell games associated with different subproblems. Putting all these pieces together, we propose a bandit online learning algorithm with the help of our bandit Blackwell approachability. We highlight that this approach essentially uses the unbiased estimators to obtain bandit-style feedback for the online learning problems of each subproblem, leading to an efficient overall bandit learning algorithm with a sublinear  $\gamma$ -regret.

### 5.1. Bandit Blackwell Sequential Games and Approachability

In the bandit online learning version of problem (1), an online algorithm can only see the value of the function at the particular point that is picked in that round. Therefore, in our transformation, multiple online Blackwell algorithms compete over a single piece of information in order to estimate the vector payoffs, where estimating the vector payoff of a Blackwell algorithm can be typically done by taking a costly “exploration” move, tailored to that algorithm.

With the goal of properly modeling this paradigm at a lower level, we propose the notion of a *bandit Blackwell sequential game*, characterized by the extended tuple  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$ . In this variant, player 1 makes an additional decision in each round: whether to *explore* or not. Only if player 1 chooses to explore in round  $t$ , do they receive the unbiased estimator  $\hat{\mathbf{p}}(\mathbf{x}_t, \mathbf{y}_t)$  whose expectation is the vector payoff for that round  $\mathbf{p}(\mathbf{x}_t, \mathbf{y}_t)$ . However, player 1 is punished by an additive cost  $D(\mathbf{p})$ . If player 1 refrains from exploration, they neither receive any feedback nor any punishment. Player 1’s new goal is to minimize the distance from the time-averaged payoff to the target set  $S$  plus their time-averaged exploration penalty.

**DEFINITION 8 (BANDIT BLACKWELL APPROACHABILITY).** A closed convex target set  $S$  is  $g(T)$ -bandit-approachable in the bandit Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$  if there exists a bandit player’s strategy such that for every player 2’s strategy, the resulting sequence of actions satisfy

$$d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) + \mathbb{E} \left[ \frac{1}{T} D(\mathbf{p}) \cdot (\# \text{ explore}) \right] \leq g(T),$$

where  $(\# \text{ explore})$  is the number of exploration rounds.

We prove the following extension of Blackwell’s approachability theorem. Interestingly, the bound in Theorem 3 in terms of the dependency on the number of rounds  $T$  is information-theoretically tight. See Section D.3 for details.

**THEOREM 3.** *A closed convex set  $S$  is  $O\left(D(\mathbf{p})^{1/3}D(\hat{\mathbf{p}})^{2/3}(\log d)^{1/3}T^{-1/3}\right)$ -bandit-approachable in the bandit Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$  if and only if  $S$  is response-satisfiable in the Blackwell game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$ . In particular, when  $S$  is response satisfiable, the online algorithm AlgBB (Algorithm 3) achieves this approachability bound in polynomial time, given access to a separation oracle for  $S$ .*

*Proof sketch of Theorem 3.* To see the only if direction of the first part of the theorem, bandit Blackwell approachability implies Blackwell approachability. Specifically, if

$$d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) + \mathbb{E} \left[ \frac{1}{T} C \cdot (\# \text{ explore}) \right] \leq O \left( D(\mathbf{p})^{1/3} D(\hat{\mathbf{p}})^{2/3} (\log d)^{1/3} T^{-1/3} \right),$$

then we must have  $d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) \leq O \left( D(\mathbf{p})^{1/3} D(\hat{\mathbf{p}})^{2/3} (\log d)^{1/3} T^{-1/3} \right)$ , and hence this  $\ell_\infty$ -distance is vanishing as  $T \rightarrow +\infty$ . This, in turn, implies that the target set  $S$  is response satisfiable (see Theorem 1). Note that while Theorem 1 is stated for a specific  $g(T)$ , the only if direction of this theorem holds for any vanishing approachability bound (Blackwell 1956).

To see the if direction and the second part of the theorem, we consider a simple algorithm that uses a (full information) Blackwell algorithm AlgB as a blackbox. We pick an algorithm AlgB that satisfies the approachability bound of Theorem 1, and can obtain this bound in polynomial time given a separation oracle for  $S$ ; see Remark 1. At the beginning of each round, our bandit algorithm plays the last suggested action by AlgB. It then explores randomly with probability  $q$  by flipping an independent coin. Based on the outcome of the coin, it either updates the state of AlgB using the unbiased payoff feedback it gets (exploration) and queries AlgB for suggesting a new action to follow, or decides not to explore with probability  $1 - q$  and refrains the state of AlgB. These steps are summarized in Algorithm 3.

As for the running time, the above algorithm will run in polynomial time given a separation oracle for  $S$  based on Remark 1. As for the approachability bound, at a high level, if we imagine that unbiased payoffs are the actual payoffs in the Blackwell game, then the expected distance of time-averaged unbiased vector payoff from  $S$  is roughly equal to the same quantity for only rounds that we explore. There are  $qT$  such rounds in expectation. Therefore, the expected distance is upper bounded by  $O(D(\hat{\mathbf{p}})(\log(d))^{1/2}q^{-1/2}T^{-1/2})$  due to the approachability of AlgB for this imaginary Blackwell sequential game (Theorem 1). Also, the algorithm gets penalized on average by  $O(D(\mathbf{p})q)$  due to exploring. Taking expectation to replace unbiased estimators with the actual payoffs and balancing the two terms in regret by setting  $q = D(\mathbf{p})^{-2/3}D(\hat{\mathbf{p}})^{2/3}(\log d)^{1/3}T^{-1/3}$  gives the final bound. See Section D.1 in the appendix for a detailed proof with a more involved argument. ■

---

**Algorithm 3:** Bandit Blackwell Online Algorithm (AlgBB)

---

**Meta Input:** Parameter  $q \in [0, 1]$ , bandit Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$ .**Input:** Number of rounds  $T$ , blackbox access to full information online algorithm AlgB for the Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$ , achieving approachability bound of Theorem 1.**Output:** Actions  $\{\mathbf{x}_t\}_{t \in [T]}$  and binary signals  $\{\pi_t\}_{t \in [T]}$ , where  $\mathbf{x}_t \in \mathcal{X}$ , and  $\pi_t \in \{\text{YES}, \text{NO}\}$  for any  $t \in [T]$ .Initialize  $\mathbf{x}_{\text{new}}$  by sending the initial query to AlgB; **for** round  $t = 1$  to  $T$  **do**    Play the action  $\mathbf{x}_t \leftarrow \mathbf{x}_{\text{new}}$ ; Set  $\pi_t$  to be YES with probability  $q$ , and NO with    probability  $1 - q$ ; **if**  $\pi_t = \text{YES}$  **then**        Obtain  $\hat{\mathbf{p}}(\mathbf{x}_t, \mathbf{y}_t)$  and send  $\hat{\mathbf{p}}(\mathbf{x}_t, \mathbf{y}_t)/q$  as feedback to AlgB; // AlgB gets a new feedback in each exploration round, i.e., round  $t$  where  $\pi_t = \text{YES}$ .    Update  $\mathbf{x}_{\text{new}}$  by querying AlgB given the actions and realized unbiased estimator vector payoffs in exploration rounds prior to round  $t + 1$ , i.e.,     $\mathbf{x}_{\text{new}} \leftarrow \text{AlgB}(\{(\mathbf{x}_\tau, \hat{\mathbf{p}}(\mathbf{x}_\tau, \mathbf{y}_\tau)) : \tau \leq t, \pi_\tau = \text{YES}\})$ ;

---

REMARK 4. Our notion of bandit Blackwell approachability and the algorithm that achieves the tight bound (Algorithm 3) bear some resemblance to the  $\epsilon$ -greedy algorithm in the classic bandit setting, where in every round of this algorithm, we decide whether or not to explore, and when we explore in a round we assume we suffer from the maximum possible regret in this round.

REMARK 5. The vanilla version of AlgBB needs to tune exploration probability  $q$  based on the horizon  $T$  to obtain the bound in Theorem 3. However, by using the standard *doubling trick* in online learning (e.g., see [Bubeck et al. \(2015\)](#)) in a blackbox fashion, one can boost Algorithm 3 to work for unknown but bounded  $T$ : the new algorithm starts with a guess for horizon (e.g.,  $T = 1$ ) and sets  $q$  according to this guess. Each time it reaches the guessed horizon, it doubles its guess, and restarts by tuning a new value for  $q$  and initializing again. The doubling trick is a well-known idea in the online learning literature that can be traced back to the classic work of [Auer et al. \(2002\)](#). We refer the reader to aforementioned work, and omit the details here for brevity.

## 5.2. Offline-to-online Transformation under Bandit Feedback

Similar to our full information offline-to-online transformation, which gave us algorithm ONLINE-IG in Section 4, we transform an offline IG algorithm to a bandit online learning algorithm by associating an instance of the bandit Blackwell sequential game to each subproblem  $i \in [N]$  of the offline algorithm. That is, we crucially rely on a reduction from the local optimization step of each subproblem in Algorithm 1 to an approachable instance of the bandit Blackwell sequential game as in Definition 8. Such a reduction is possible if the offline algorithm is *Bandit Blackwell reducible*; see the following definition.

DEFINITION 9 (BANDIT BLACKWELL REDUCIBILITY). An instance OFFLINE-IG( $\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta$ ) of Algorithm 1 is bandit Blackwell reducible if there is an instance  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$  of bandit Blackwell sequential game (Section 5.1) and an *exploration sampling device*  $\text{ExpS}: \Theta \times \mathcal{D} \rightarrow \Delta(\mathbb{R}^{d_{\text{payoff}}} \times \mathcal{C})$ , such that:

1. OFFLINE-IG( $\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta$ ) is Blackwell reducible as in Definition 7, using the Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  (with biaffine  $\mathbf{p}$ ) and the synthetic Blackwell adversary function AdvB.
2. If  $\mathbf{y} = \text{AdvB}(\mathbf{z}, f)$  for some  $f \in \mathcal{F}, \mathbf{z} \in \mathcal{D}$ , then  $\hat{\mathbf{p}}(\boldsymbol{\theta}, \mathbf{y}) = f(\mathbf{z}_{\text{exp}})\mathbf{w}_{\text{exp}}$  for all  $\boldsymbol{\theta} \in \Theta$ , where  $(\mathbf{w}_{\text{exp}}, \mathbf{z}_{\text{exp}}) \sim \text{ExpS}(\boldsymbol{\theta}, \mathbf{z})$ . Otherwise,  $\hat{\mathbf{p}}(\boldsymbol{\theta}, \mathbf{y}) = \mathbf{p}(\boldsymbol{\theta}, \mathbf{y})$ .
3. The above  $\hat{\mathbf{p}}$  is an unbiased estimator for the actual vector payoff, i.e., for all  $\boldsymbol{\theta} \in \Theta, \mathbf{y} \in \mathcal{Y}$ :  $\mathbb{E}[\hat{\mathbf{p}}(\boldsymbol{\theta}, \mathbf{y})] = \mathbf{p}(\boldsymbol{\theta}, \mathbf{y})$ .
4. The exploration sampling device  $\text{ExpS}(\boldsymbol{\theta}, \mathbf{z})$  returns its samples  $(\mathbf{w}_{\text{exp}}, \mathbf{z}_{\text{exp}})$  in polynomial time.

To better understand the bandit Blackwell reducibility, we revisit our running example.

EXAMPLE 1 (CONTINUED). The greedy algorithm of Nemhauser et al. (1978) is also bandit Blackwell reducible. As stated in Section 4, this algorithm is Blackwell reducible. Recall that in this example, the biaffine Blackwell payoff is  $\mathbf{p}(\boldsymbol{\theta}, \mathbf{y}) = \boldsymbol{\theta}^T \mathbf{y} \mathbf{1}_n - \mathbf{y}$ , where  $\mathbf{1}_n$  is all ones  $n$ -dimensional vector. We will construct an exploration sampling device ExpS that returns  $(\mathbf{w}_{\text{exp}}, \mathbf{z}_{\text{exp}})$  such that if  $\forall \boldsymbol{\theta} \in \Theta$ , we have  $\mathbf{y} = \text{AdvB}(\mathbf{z}, f)$  for some  $f \in \mathcal{F}, \mathbf{z} \in \mathcal{D}$ , we set  $\hat{\mathbf{p}}(\boldsymbol{\theta}, \mathbf{y}) = f(\mathbf{z}_{\text{exp}})\mathbf{w}_{\text{exp}}$  and we must have  $\mathbb{E}[\hat{\mathbf{p}}(\boldsymbol{\theta}, \mathbf{y})] = \mathbf{p}(\boldsymbol{\theta}, \mathbf{y})$ . The exploration sampling device ExpS works as follows. Given a point  $\mathbf{z} \in \mathcal{C}$  (which represents a set of elements) and parameter  $\boldsymbol{\theta} \in \Theta$ , it draws  $j \sim \text{Uniform}\{1, \dots, n\}$  and returns (i)  $\mathbf{w}_{\text{exp}} = n(\boldsymbol{\theta}_j \mathbf{1}_n - \mathbf{e}_j)$ , (ii)  $\mathbf{z}_{\text{exp}} = \mathbf{z} \cup \{j\}$ . Now,  $\hat{\mathbf{p}}$  is an unbiased estimator of  $\mathbf{p}$ , because:

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{p}}(\boldsymbol{\theta}, \text{AdvB}(\mathbf{z}, f))] &= \mathbb{E}[f(\mathbf{z}_{\text{exp}})\mathbf{w}_{\text{exp}}] \\
&= \mathbb{E}[n(\boldsymbol{\theta}_j f(\mathbf{z} \cup \{j\})\mathbf{1}_n - f(\mathbf{z} \cup \{j\})\mathbf{e}_j)] \\
&= \mathbf{1}_n \sum_{j \in [n]} \boldsymbol{\theta}_j f(\mathbf{z} \cup \{j\}) - [f(\mathbf{z} \cup \{1\}), \dots, f(\mathbf{z} \cup \{n\})]^T \\
&= \mathbf{1}_n \sum_{j \in [n]} \boldsymbol{\theta}_j (f(\mathbf{z} \cup \{j\}) - f(\mathbf{z})) - [f(\mathbf{z} \cup \{1\}), \dots, f(\mathbf{z} \cup \{n\})]^T + f(\mathbf{z})\mathbf{1}_n \\
&= \boldsymbol{\theta}^T \mathbf{y} \mathbf{1}_n - \mathbf{y} = \mathbf{p}(\boldsymbol{\theta}, \text{AdvB}(\mathbf{z}, f)),
\end{aligned}$$

where  $\mathbf{y} \triangleq [f(\mathbf{z} \cup \{j\}) - f(\mathbf{z})]_{j=1,2,\dots,n} = \text{AdvB}(\mathbf{z}, f)$ . Here, the fourth equation holds because  $\sum_{j \in [n]} \boldsymbol{\theta}_j = 1$ . Observe that the exploration sampling device ExpS has an intuitive interpretation, at every round, it randomly picks one of the elements  $j \in [n]$ , and evaluates the marginal benefit of adding element  $j$  to  $\mathbf{z}$ .

**5.2.1. Overview of the Algorithm** When the offline algorithm (Algorithm 1) is bandit Blackwell reducible (Definition 9), we can employ a similar offline-to-online transformation mentioned in Section 4. However, instead of associating an instance of the Blackwell game to each subproblem, we associate an instance of the bandit Blackwell game. To obtain unbiased estimators for the vector payoffs of these bandit Blackwell instances, we rely on the exploration sampling devices that are promised by Definition 9. This sampling device allows us to strike a balance between exploration and exploitation in all of the online bandit Blackwell games. We formalize this transformation of the offline algorithm to an online bandit algorithm called BANDIT-IG in Algorithm 4.

Suppose that the offline algorithm OFFLINE-IG( $\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta$ ) is given. For the particular bandit Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$  coming from Definition 9, we use AlgBB (Algorithm 3) to determine the strategy of player 1. Such an online bandit Blackwell algorithm as player 1 ensures that the distance between the average vector payoff  $\frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t)$  and set  $S$  plus the exploration penalty goes to zero with rate  $g(T) = O(D(\mathbf{p})^{1/3} D(\hat{\mathbf{p}})^{2/3} (\log d)^{1/3} T^{-1/3})$ ; see Theorem 3.

We dedicate a copy of the above algorithm AlgBB<sup>(i)</sup> to each subproblem  $i \in [N]$ . We query algorithms AlgBB<sup>(i)</sup> in the increasing order of their index  $i$ . Consider the online bandit Blackwell algorithm AlgBB<sup>(i)</sup>, and assume that in round  $t$ , we query this algorithm. The algorithm returns two outputs: the update parameter  $\boldsymbol{\theta}_t^{(i)}$  and a binary signal  $\pi_t^{(i)} \in \{\text{YES}, \text{NO}\}$ . If  $\pi_t^{(i)} = \text{YES}$ , the algorithm explores: it samples  $(\mathbf{w}_{t,\text{exp}}^{(i)}, \mathbf{z}_{t,\text{exp}}^{(i)})$  from the exploration sampling device  $\text{ExpS}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)})$ . Note that the exploration sampling device uses the update parameter  $\boldsymbol{\theta}_t^{(i)}$  and the point returned by the previous subproblem  $\mathbf{z}_t^{(i-1)}$ . This indeed allows the subproblems to communicate with each other during exploration. The algorithm then plays  $\mathbf{z}_t = \mathbf{z}_{t,\text{exp}}^{(i)}$  and provides the payoff vector feedback  $\hat{\mathbf{p}}_t^{(i)} = f_t(\mathbf{z}_t) \mathbf{w}_{t,\text{exp}}^{(i)}$  to AlgBB<sup>(i)</sup>. This feedback is only used by the online bandit Blackwell algorithm AlgBB<sup>(i)</sup>, not the rest of  $N - 1$  bandit Blackwell algorithms. We highlight that if AlgBB<sup>(i)</sup> decides to explore in round  $t$ , the rest of bandit Blackwell algorithms will not be queried. Finally, if  $\pi_t^{(i)} = \text{NO}$ , the algorithm exploits: it returns point  $\mathbf{z}_t^{(i)} = \text{LOCAL-UPDATE}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)})$ . Again observe that during exploitation, subproblem  $i$  also communicates with subproblem  $i - 1$  through using  $\mathbf{z}_t^{(i-1)}$ .

**5.2.2. Regret Analysis** Theorem 4 bounds the regret of the BANDIT-IG algorithm. The proof is deferred to Section D.2 in the appendix.

**THEOREM 4 (Bandit information offline-to-online transformation).** *Suppose that an instance of OFFLINE-IG( $\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta$ ) for the offline problem (1) satisfies the following properties:*

- *It is an extended  $(\gamma, \delta)$ -robust approximation for  $\gamma \in (0, 1)$  and  $\delta > 0$ , as in Definition 6.*
- *It is bandit Blackwell reducible; that is, we can define the bandit Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$  and exploration sampling device  $\text{ExpS} : \Theta \times \mathcal{D} \rightarrow \Delta(\mathbb{R}^{d_{\text{payoff}}} \times \mathcal{C})$  that satisfy the conditions in Definition 9.*

**Algorithm 4:** Bandit Online Learning Meta-algorithm (BANDIT-IG)

**Meta Input:** Feasible region  $\mathcal{C}$ , function space  $\mathcal{F}$ , defined over domain  $\mathcal{D}$ , parameter space

$$\Theta \subseteq \mathbb{R}^{d_{\text{param}}}.$$

**Offline algorithm and reduction gadgets:** An instance OFFLINE-IG( $\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta$ ) of Algorithm 1; this algorithm is bandit Blackwell reducible as in Definition 9, using the bandit Blackwell instance  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$  and exploration sampling device

$$\text{ExpS} : \Theta \times \mathcal{D} \times \rightarrow \Delta(\mathbb{R}^{d_{\text{payoff}}} \times \mathcal{C}).$$

**Input:** Number of rounds  $T$ ; access to a bandit Blackwell online algorithm AlgBB.

**Output:** Points  $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_T \in \mathcal{C}$ .

Initialize  $N$  parallel instances  $\{\text{AlgBB}^{(i)}\}_{i=1}^N$  of the online algorithm AlgBB; **for** round  $t = 1$

**to**  $T$  **do**

Initialize  $\mathbf{z}_t^{(0)} \in \mathcal{C}$ ; **for** subproblem  $i = 1$  **to**  $N$  **do**

Choose the update parameter  $\boldsymbol{\theta}_t^{(i)} \in \Theta$  and exploration signal  $\pi_t^{(i)} \in \{\text{YES}, \text{NO}\}$  by querying online algorithm AlgBB<sup>(i)</sup> given the update parameters and vector payoffs  $\hat{\mathbf{p}}$  of exploration rounds prior to round  $t$  in the bandit Blackwell sequential game of subproblem  $i$ , that is

$$\left( \boldsymbol{\theta}_t^{(i)}, \pi_t^{(i)} \right) \leftarrow \text{AlgBB}^{(i)} \left( \boldsymbol{\theta}_1^{(i)}, \dots, \boldsymbol{\theta}_{t-1}^{(i)}, \{ \hat{\mathbf{p}}(\boldsymbol{\theta}_\tau^{(i)}, \mathbf{y}_\tau^{(i)}) \}_{\tau \leq t-1; \pi_\tau^{(i)} = \text{YES}} \right); \text{ if } \pi_t^{(i)} = \text{YES},$$

**then**

Sample  $(\mathbf{w}_{t,\text{exp}}^{(i)}, \mathbf{z}_{t,\text{exp}}^{(i)})$  from the exploration sampling device  $\text{ExpS}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)})$ ;

Play the exploration point  $\mathbf{z}_t \leftarrow \mathbf{z}_{t,\text{exp}}^{(i)}$ ;

$\langle$  Bandit information feedback: observe  $f_t(\mathbf{z}_t) \rangle$ ;

Give payoff vector feedback  $\hat{\mathbf{p}}_t^{(i)} = f_t(\mathbf{z}_t) \cdot \mathbf{w}_{t,\text{exp}}^{(i)}$  to AlgBB<sup>(i)</sup>; Skip immediately to the beginning of the next round  $t + 1$ ;

Set  $\mathbf{z}_t^{(i)} \leftarrow \text{LOCAL-UPDATE}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)})$ ;

Play the final point  $\mathbf{z}_t \leftarrow \mathbf{z}_t^{(N)}$  and receive bandit feedback  $f_t(\mathbf{z}_t)$ , and ignore it.

Consider the bandit-information adversarial online learning version of problem (1), and let AlgBB be a polynomial time bandit Blackwell algorithm for  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$  as in Theorem 3. Then, for this online problem, BANDIT-IG( $\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta, \text{AlgBB}$ ) runs in polynomial time and satisfies the following  $\gamma$ -regret bound:

$$\gamma\text{-regret}(\text{BANDIT-IG}(\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta, \text{AlgBB})) \leq O\left(D(\mathbf{p})^{1/3} D(\hat{\mathbf{p}})^{2/3} N \delta (\log(d_{\text{payoff}}))^{1/3} T^{2/3}\right),$$

where  $N$  is the number of subproblems and  $d_{\text{payoff}}$  is the dimension of vector payoffs.

We finish this section by wrapping up our running example (Example 1) and mentioning the bandit regret bound we get as a direct corollary of Theorem 4.

EXAMPLE 1 (FINISHED). The greedy algorithm in Nemhauser et al. (1978) satisfies  $(1 - \frac{1}{e}, 1)$ -robust approximation and is bandit Blackwell reducible. It has  $N = k$  subproblems and  $\ell_\infty$  diameter of  $\hat{\mathbf{p}}$  is  $D(\hat{\mathbf{p}}) = O(n)$ . Therefore, by invoking Algorithm 4 given any bandit Blackwell algorithm satisfying the approachability bound in Theorem 3, we obtain the following bandit regret bound:

$$\left(1 - \frac{1}{e}\right)\text{-regret}(\text{Algorithm 4}) \leq O(kn^{2/3}(\log n)^{1/3}T^{2/3}),$$

which in turn, by noting that  $k$  can be as large as  $n$ , gives us an immediate improvement over regret bound of  $O(k^2(n \log n)^{1/3}T^{2/3}(\log T)^2)$  in Streeter and Golovin (2007, 2008).

## 6. Applications to Revenue Management and Combinatorial Optimization

We have already showed how to fit monotone submodular maximization into our framework through Example 1. In this section, we apply our framework to two other selected problems: product ranking through sequential submodular maximization and personalized reserve price optimization in second-price auction. (See Section E in the appendix to see how to apply our framework to the problem of maximizing non-monotone continuous weak-DR submodular functions, in which obtaining any sub-linear approximate regret has been an open problem for a while. See also Section F to see how to apply our framework to the problem of maximizing monotone continuous strong-DR submodular functions subject to downward closed bounded convex sets.) Our framework results in improved/new regret bounds in all mentioned applications for both full-information and bandit settings.

### 6.1. Application to Product Ranking and Sequential Submodular Maximization

*Problem definition.* In the PRODUCT RANKING PROBLEM, a platform aims to characterize a ranking of  $n$  items, where a ranking is a permutation  $\boldsymbol{\pi}$  over the items. Here, items on positions with lower indices have more visibility. The goal of the platform is to maximize its user engagement (also known as market share), which is the probability that a consumer does not leave the platform without taking a desired action. This action can be a click, purchase, or even installing an application.

For the sake of presentation, assume that the desired action is clicking on an item. We consider the model proposed by Asadpour et al. (2020), which is inspired by an earlier model proposed in Ferreira et al. (2021). In this model, a consumer  $u$  is characterized by a patience level  $\theta_u$  together with a monotone non-decreasing submodular set function  $\kappa_u : 2^{[n]} \rightarrow [0, 1]$ . A consumer of type  $(\theta_u, \kappa_u)$ , when offered a ranked list of products  $\boldsymbol{\pi} = ([\boldsymbol{\pi}]_1, [\boldsymbol{\pi}]_2, \dots, [\boldsymbol{\pi}]_n)$ , inspects the first  $\theta_u$  products and clicks with probability  $\kappa_u(\{[\boldsymbol{\pi}]_1, \dots, [\boldsymbol{\pi}]_{\theta_u}\})$ . The platform knows the distribution  $\mathcal{G}$  from which  $u$  is selected. The goal is to pick a permutation  $\boldsymbol{\pi}$  maximizing the probability of click

$$\mathbb{E}_{u \sim \mathcal{G}} [\kappa_u(\{[\boldsymbol{\pi}]_1, \dots, [\boldsymbol{\pi}]_{\theta_u}\})].$$

For a wide range of choice models in the literature, the probability of a purchase from an offered set  $S$  can be described using a monotone submodular function  $\kappa_u$ . This includes multinomial logit, nested logit, and paired combinatorial logit models. See [Kök et al. \(2008\)](#) for details on these models.

*Product ranking problem as sequential submodular maximization.* A slight reformulation of the above model casts the product ranking problem as a special case of a class of optimization problems over permutations called sequential submodular maximization ([Asadpour et al. 2020](#)). We define the sequential submodular maximization problem as follows.<sup>15</sup> Given a sequence of monotone submodular set functions  $\{f_1(\cdot), \dots, f_n(\cdot)\}$ , and a sequence of non-negative weights  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)$ , we aim to find a ranking  $\pi$  that maximizes

$$\sum_{i=1}^n \lambda_i f_i(\{[\boldsymbol{\pi}]_1, \dots, [\boldsymbol{\pi}]_i\}),$$

where  $[\boldsymbol{\pi}]_i$  denotes the item on the  $i^{\text{th}}$  position of ranking  $\pi$ . In the aforementioned choice model, for all  $i \in [n]$ , we have  $f_i(S) \triangleq \mathbb{E}_{u \sim \mathcal{G}} [\kappa_u(S) | \theta_u = i]$ , representing the probability of clicks functions, and  $\lambda_i \triangleq \mathbb{P}_{u \sim \mathcal{G}} (\theta_u = i)$ , representing the probability that a consumer has patience level  $i$ . The probability that a consumer clicks on at least one product when offered a ranked set of products  $\boldsymbol{\pi}$  is then

$$f(\boldsymbol{\pi}) \triangleq \lambda_1 f_1(\{[\boldsymbol{\pi}]_1\}) + \lambda_2 f_2(\{[\boldsymbol{\pi}]_1, [\boldsymbol{\pi}]_2\}) + \dots + \lambda_n f_n(\{[\boldsymbol{\pi}]_1, \dots, [\boldsymbol{\pi}]_n\}),$$

where  $f_i$ 's are monotone submodular functions and  $\lambda_i$ 's are non-negative. To simplify the analysis, notice that while  $f$  is a function of a set of ranked/ordered items,  $f_i$  is a function of a set that has at most  $i$  items for each  $i \in [n]$ .

*Online problem.* In the offline setting, the platform knows  $\mathcal{G}$ , which translates to knowing the probability of clicks  $\{f_1(\cdot), \dots, f_n(\cdot)\}$  and the probability distribution of the patience level  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)$ . We study the online user-engagement-maximization ranking problem where on every round  $t$ , a distribution over patience level  $\boldsymbol{\lambda}_t$  and the expected probability of click function  $f_t$ , which is made of  $\{f_{t,1}(\cdot), \dots, f_{t,n}\}$ , are chosen adversarially. The platform, whose goal is to maximize its user-engagement, chooses a ranking  $\pi_t$  without observing  $\boldsymbol{\lambda}_t$  and  $f_t$ . After choosing the ranking, the platform observes the function  $f_t$  in the full-information setting. In the bandit setting, the platform *only* observes whether or not the consumer clicks on at least one item, but not which item was clicked on. To the best of our knowledge, the online adversarial version of this problem has not been studied before, neither under full information nor bandit setting.

<sup>15</sup> Our notion of sequential submodular functions is different from the notion of sequential submodular functions studied in [Tschitschek et al. \(2017\)](#), [Mitrovic et al. \(2018\)](#). However, under all these notions, the goal is to return a sequence (permutation), rather than a set.

*Offline algorithm.* In this paper, we focus on a greedy algorithm that achieves  $\frac{1}{2}$ -approximation, and transform it into an online adversarial learning algorithm. Our offline algorithm is presented in Algorithm 5. The input to this algorithm is a sequential submodular function  $f : \Pi \rightarrow [0, 1]$ , where  $\Pi$  is the set of ranking permutations of  $n$  items, i.e.,  $\Pi = \{0, 1, \dots, n\}^n$ . In this case, having  $[\boldsymbol{\pi}]_i = 0$  represents putting no item at position  $i$  and multiple positions can display the same item for simplicity. In this problem, both the domain and the feasible region are  $\mathcal{D} = \mathcal{C} = \Pi$ . Let  $S_i$  denote the set of subsets of  $[n]$  that consist of at most  $i$  items, i.e.,  $S_i = \{S \subseteq [n] \mid |S| \leq i\}$ . We have  $f(\boldsymbol{\pi}) = \sum_{j=1}^n \lambda_j f_j(\{[\boldsymbol{\pi}]_1, \dots, [\boldsymbol{\pi}]_j\})$  where  $f_i$  is a monotone submodular function that takes an element of  $S_i$  as input and returns a probability in  $[0, 1]$ , i.e.,  $f_i : S_i \rightarrow [0, 1]$ . Algorithm 5, taken from Asadpour et al. (2020) and Ferreira et al. (2021), is a greedy algorithm with  $n$  subproblems, where each subproblem corresponds to a position on the ranking. The algorithm starts filling up the positions from the top and for each position  $i$ , it chooses the item that has the highest marginal probability of click. The update  $\boldsymbol{\pi}^{(i)} \leftarrow \boldsymbol{\pi}^{(i-1)} + z^{(i)} \mathbf{e}_i$  represents the action of adding item  $z^{(i)}$  to position  $i$ .

---

**Algorithm 5:** Greedy for Sequential Submodular Maximization (Asadpour et al. 2020)

---

**Input:** A sequential submodular function  $f$ , which can be represented using a sequence of monotone submodular functions  $\{f_i(\cdot)\}_{i \in [n]}$  and a sequence of non-negative weights  $\boldsymbol{\lambda}$ .

**Output:** Ranking  $\boldsymbol{\pi} \in \Pi$ .

Set initial ranking  $\boldsymbol{\pi}^{(0)} \leftarrow \mathbf{0}_n$ .

**for** position  $i = 1, 2, \dots, n$  **do**

Local Optimization Step

Choose  $z^{(i)} \in \arg \max_{z \in [n]} \sum_{j=i}^n \lambda_j f_j(\{[\boldsymbol{\pi}^{(i-1)}]_1, \dots, [\boldsymbol{\pi}^{(i-1)}]_{i-1}, z\}) - \sum_{j=i}^n \lambda_j f_j(\{[\boldsymbol{\pi}^{(i-1)}]_1, \dots, [\boldsymbol{\pi}^{(i-1)}]_{i-1}\})$ .

Local Update Step

Set  $\boldsymbol{\pi}^{(i)} \leftarrow \boldsymbol{\pi}^{(i-1)} + z^{(i)} \mathbf{e}_i$ .

**return**  $\boldsymbol{\pi} \leftarrow \boldsymbol{\pi}^{(n)}$ .

---

We cast Algorithm 5 as an instance of OFFLINE-IG (Algorithm 1). The parameter space is  $\Theta = \Delta([n])$  and  $d_{\text{param}} = n$ . Moreover, in subproblem  $i$  the algorithm picks the distribution  $\boldsymbol{\theta}^{(i)}$  over items so that the resulting vector payoff lands in set  $S$ . In this language, set  $S$  is the  $n$ -dimensional positive orthant and the vector payoff function is:

$$\forall j \in [n]: \quad [\text{PAYOFF}(\boldsymbol{\theta}^{(i)}, \boldsymbol{\pi}^{(i-1)}, f)]_j = \boldsymbol{\theta}^T \mathbf{y}^{(i)} - [\mathbf{y}^{(i)}]_j \quad ,$$

where

$$\begin{aligned} \mathbf{y}^{(i)} &\triangleq \left[ \sum_{a=i}^n \lambda_a f_a(\{[\boldsymbol{\pi}^{(i-1)}]_1, \dots, [\boldsymbol{\pi}^{(i-1)}]_{i-1}, j\}) - \sum_{a=i}^n \lambda_a f_a(\{[\boldsymbol{\pi}^{(i-1)}]_1, \dots, [\boldsymbol{\pi}^{(i-1)}]_{i-1}\}) \right]_{j \in [n]} \\ &= [f(\boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_i) - f(\boldsymbol{\pi}^{(i-1)})]_{j \in [n]}, \end{aligned}$$

is the marginal objective value of putting item  $j$  on the  $i^{\text{th}}$  position. Note that any  $\boldsymbol{\theta}^{(i)}$  for which the vector payoff is in  $S$  is indeed a distribution over items  $z^{(i)}$  such that  $z^{(i)} \in \arg \max_{z \in [n]} \sum_{j=i}^n \lambda_j f_j(\{[\boldsymbol{\pi}^{(i-1)}]_1, \dots, [\boldsymbol{\pi}^{(i-1)}]_{i-1}, z\}) - \sum_{j=i}^n \lambda_j f_j(\{[\boldsymbol{\pi}^{(i-1)}]_1, \dots, [\boldsymbol{\pi}^{(i-1)}]_{i-1}\})$ .

**THEOREM 5 (Online learning for sequential submodular maximization).** *Let  $n$  be the number of items. For the problem of maximizing sequential submodular functions in the online full information setting, there exists a learning algorithm that obtains  $O(n\sqrt{T \log n})$   $\frac{1}{2}$ -regret, where  $T$  is the number of rounds. Furthermore, in the online bandit setting, there exists a learning algorithm that obtains  $O(n^{5/3}(\log n)^{1/3} T^{2/3})$   $\frac{1}{2}$ -regret, where  $T$  is the number of rounds. For this problem, the benchmark in the regret bounds is  $\frac{1}{2} \max_{\boldsymbol{\pi} \in \Pi} \sum_{t=1}^T f_t(\boldsymbol{\pi})$ .*

Theorem 5 and the following corollary are proved in Section G using our offline-to-online transformations presented in Sections 4 and 5.

**COROLLARY 1 (Online learning for product ranking).** *Let  $n$  be the number of items. For the problem of product ranking optimization to maximize user engagement using the model from Asadpour et al. (2020), there exists a learning algorithm that obtains  $O(n\sqrt{T \log n})$   $\frac{1}{2}$ -regret in the full-information setting and  $O(n^{5/3}(\log n)^{1/3} T^{2/3})$   $\frac{1}{2}$ -regret in the bandit setting, where  $T$  is the number of consumers.<sup>16</sup>*

In Section A.1, we numerically evaluate of our online algorithms for the product ranking problem under both full information and bandit settings.

## 6.2. Application to Maximizing Multiple Reserves in Second Price Auction

*Problem definition.* In the MAXIMIZING MULTIPLE RESERVES (MMR) problem (Roughgarden and Wang 2019, Derakhshan et al. 2019), a seller wants to sell one item to  $n$  bidders to maximize her revenue. Each bidder  $i$  has a private value  $v_i$  for the item. The seller runs a second-price auction with personalized reserves  $\mathbf{r}$ ; the winner is the bidder with the highest bid/valuation among the bidders whose bids exceed their reserve prices.<sup>17</sup> The winner pays the minimum bid she needed to win, which is the maximum between their reserve price and the second-highest bid that cleared its reserve price.

<sup>16</sup> As an implication, the same problem for the consumer choice model from Ferreira et al. (2021) also has a learning algorithm that obtains  $O(n\sqrt{T \log n})$   $\frac{1}{2}$ -regret in the full information setting and  $O(n^{5/3}(\log n)^{1/3} T^{2/3})$   $\frac{1}{2}$ -regret in the bandit setting. For both problems, the benchmark in the regret bounds is  $\frac{1}{2} \max_{\boldsymbol{\pi} \in \Pi} \sum_{t=1}^T f_t(\boldsymbol{\pi})$ .

<sup>17</sup> Since second price auctions are truthful, we will use bids and valuations interchangeably.

*Online problem.* We are interested in the seller’s problem in the online full information and bandit settings. In both settings, each round  $t \in [T]$  involves the seller choosing a set of reserves  $\mathbf{r}_t$  and the adversary choosing a valuation profile  $\mathbf{v}_t$ . In the online full information setting, the seller observes the valuation profile and gets credit for the resulting revenue. In the online bandit setting, the seller observes just the resulting revenue and does not observe the bidders’ valuations or even the identity of the winner. The seller’s goal is to minimize the difference between his average revenue and the best average revenue in hindsight for a fixed set of reserves  $\mathbf{r}^*$ . To the best of our knowledge, the bandit setting has not been studied in the literature. However, the full information setting of this problem is studied by [Roughgarden and Wang \(2019\)](#), where they present a learning algorithm with  $O(n\sqrt{T\log T})$   $\frac{1}{2}$ -regret, which we improve upon.

*Offline non-batch vs batch problem.* We start with formulating the offline non-batch problem. Let  $\mathcal{R} = \{\rho_1, \rho_2, \dots, \rho_m\}$  be the set of feasible reserve prices, where  $|\mathcal{R}| = m$ , and they are sorted:  $0 = \rho_1 < \rho_2 < \dots < \rho_m$ . For the offline (non-batch) problem, let  $f : \mathcal{R}^n \times [0, 1]^n \rightarrow [0, 1]$  be the seller’s revenue function:  $f(\mathbf{r}, \mathbf{v}) = \max\{[\mathbf{v}]_{\hat{j}}, [\mathbf{r}]_{j^*}\}$  for some  $\mathbf{v}$ . Here,  $j^*$  and  $\hat{j}$  are the highest and second-highest bidders among those who cleared their bids, with ties broken arbitrarily:  $j^* \in \arg \max_{j \in [n]: [\mathbf{v}]_j \geq [\mathbf{r}]_j} \{[\mathbf{v}]_j\}$  and  $\hat{j} \in \arg \max_{j \in [n]: [\mathbf{v}]_j \geq [\mathbf{r}]_j, j \neq j^*} \{[\mathbf{v}]_j\}$ . If no bidder clears their reserve, then we say  $[\mathbf{r}]_{j^*}$  and  $[\mathbf{v}]_{\hat{j}}$  are both zero. Similarly, if only one bidder clears their reserve, we say  $[\mathbf{v}]_{\hat{j}}$  is zero. Moreover,  $\mathcal{F}$  is the space of all such revenue functions:  $\mathcal{F} = \{f \mid \exists \mathbf{v} \in [0, 1]^n \text{ such that } f(\mathbf{r}, \mathbf{v}) = \max\{[\mathbf{v}]_{\hat{j}}, [\mathbf{r}]_{j^*}\}\}$ . In the offline (non-batch) problem, the goal is to solve  $\max_{\mathbf{r} \in \mathcal{R}^n} f(\mathbf{r}, \mathbf{v})$  for an input valuation profile  $\mathbf{v} \in [0, 1]^n$ . In the optimization problem, the domain  $\mathcal{D}$  we consider and the feasible region  $\mathcal{C}$  are both  $\mathcal{R}^n$ .

The aforementioned offline problem can be solved efficiently. Note that in the offline problem, the seller who has access to the valuations of the bidders in one auction needs to optimize personalized reserve prices. It is then obvious that in the offline setting, the best action is to set reserve prices of all the bidders to zero, except the bidder with the highest bid; for this bidder, his reserve price is set to his valuation. Then, one may wonder why for the online version of this offline (non-batch) problem, which is not even NP-hard, we characterize  $\frac{1}{2}$ -regret, rather than 1-regret. The reason is that [Roughgarden and Wang \(2019\)](#) show that the full information online setting is at least as hard as the offline batch problem, which is APX-hard. In the offline batch problem, the seller has access to the valuation profiles in  $m$  auctions and would like to determine a single vector of reserve prices  $\mathbf{r}$  that maximizes revenue across all the  $m$  auctions. Considering the hardness of the offline batch problem, to solve the offline (non-batch) problem, we use a slight variation of the algorithm of [Roughgarden and Wang \(2019\)](#). This variation is stated in Algorithm 6, which obtains a  $\frac{1}{2}$  fraction of the optimal revenue similar to the original algorithm. See Section H.1 for a discussion on the major differences between the two.

**Algorithm 6:** Greedy Algorithm for Discretized MMR (Roughgarden and Wang 2019)**Input:** Valuation profile  $\mathbf{v}$ .**Output:** Reserve prices  $\mathbf{r} \in \mathcal{R}^n$ .Set initial reserves  $\mathbf{r}^{(0)} \leftarrow \mathbf{0}_n$ .**for** bidder  $i = 1, 2, \dots, n$  **do**

Define *revenue-from-reserves* function  $q^{(i)} : \mathcal{R} \rightarrow [0, 1]$  as  $q^{(i)}(r)$  equals  $r$  if  $i$  has the highest valuation (ties broken arbitrarily) and  $r \in [[\mathbf{v}]_{i'}, [\mathbf{v}]_i]$  where  $i'$  has the second-highest valuation, and 0 otherwise.

Local Optimization Step

Choose  $z^{(i)} \in \arg \max_{r \in \mathcal{R}} q^{(i)}(r)$ . // In this case  $\theta^{(i)} \in \Delta(\mathcal{R})$  is the distribution that always returns  $z^{(i)}$ .

Local Update Step

Set  $\mathbf{r}^{(i)} \leftarrow \mathbf{r}^{(i-1)} + z^{(i)} \mathbf{e}_i$ .

**return**  $\mathbf{r} \sim \text{Uniform}\{\mathbf{0}_n, \mathbf{r}^{(n)}\}$ .

*Offline algorithm.* We now briefly discuss Algorithm 6 and show how to cast it as an instance of OFFLINE-IG (Algorithm 1). This greedy algorithm has  $n$  subproblems, where in each subproblem, reserve price of a bidder  $i$  is set using our *revenue-from-reserves* function  $q$ . At the end, the algorithm randomly returns either the all-zeros reserve vector  $\mathbf{0}_n$  or the crafted reserve vector denoted by  $\mathbf{r}^{(n)}$ , where the former yields revenue equal to the second-highest valuation and the latter yields revenue of at least  $q^{(j^*)}(z^{(j^*)})$ ; see the definition of  $q(\cdot)$  in the algorithm. By definition of the revenue function, the optimal reserves obtain their revenue via one of these two cases, i.e.,

$$f(\mathbf{r}^*, \mathbf{v}) \leq \max\{[\mathbf{v}]_j, q^{(j^*)}([\mathbf{r}]_{j^*})\} \leq [\mathbf{v}]_j + q^{(j^*)}([\mathbf{r}]_{j^*}) \leq f(\mathbf{0}_n, \mathbf{v}) + f(\mathbf{r}^{(n)}, \mathbf{v}) = 2\mathbb{E}[f(\mathbf{r}, \mathbf{v})],$$

where the expectation is taken with respect to the randomness in the algorithm. This implies that our algorithm is indeed a  $\frac{1}{2}$ -approximation. Stated in the language of our Algorithm 1, our local updates manage to guarantee that  $\text{PAYOFF}(\theta^{(i)}, \mathbf{r}^{(i-1)}, \mathbf{v})$  is in the positive orthant, where the (asymmetric) vector payoff function  $\text{PAYOFF}$  returns an  $m$ -dimensional point whose  $j^{\text{th}}$  coordinate value is the expected difference between the expected value of picking a reserve according to  $\theta^{(i)}$  and that of picking  $\rho_j$ :  $[\text{PAYOFF}(\theta^{(i)}, \mathbf{r}^{(i-1)}, \mathbf{v})]_j \triangleq \mathbb{E}_{z' \sim \theta^{(i)}} [q^{(i)}(z') - q^{(i)}(\rho_j)]$ .<sup>18</sup> The following theorem shows that using our framework, the greedy Algorithm 6 can be transformed to polynomial-time online learning algorithms under both full information and bandit feedback structures. See Section A.2 for numerical evaluation of our online algorithms for reserve price optimization under both full information and bandit settings.

<sup>18</sup> Note that here,  $\text{PAYOFF}(\theta^{(i)}, \mathbf{r}^{(i-1)}, \mathbf{v})$  is not a function of  $\mathbf{r}^{(i-1)}$ .

**THEOREM 6 (Online learning for maximizing multiple reserves).** *Let  $\mathcal{R} = \{\rho_1, \dots, \rho_m\}$  be the set of possible reserve prices and  $n$  be the number of bidders. Assume that the maximum valuation is normalized to one. Then, for the problem of maximizing personalized reserve prices in the online full information setting, there exists a learning algorithm that obtains  $O\left(nT^{1/2} \log^{1/2} m\right)$   $\frac{1}{2}$ -regret, where  $T$  is the number of auctions. Furthermore, in the online bandit setting, there exists a learning algorithm that obtains  $O\left(nm^{2/3}T^{2/3} \log^{1/3} m\right)$   $\frac{1}{2}$ -regret, where  $T$  is the number of auctions. Here, the benchmark in the regret bounds is  $\frac{1}{2} \max_{\mathbf{r} \in \mathcal{R}^n} \sum_{t=1}^T f(\mathbf{r}, \mathbf{v}_t)$ .*

The proof of Theorem 6 is presented in Section H.2. The following corollary considers a stronger benchmark than the one we considered earlier. This benchmark allows the reserve prices to be any number in  $[0, 1]^n$ , i.e., the regret is computed against  $\frac{1}{2} \max_{\mathbf{r} \in [0, 1]^n} \sum_{t=1}^T f(\mathbf{r}, \mathbf{v}_t)$ , rather than  $\frac{1}{2} \max_{\mathbf{r} \in \mathcal{R}^n} \sum_{t=1}^T f(\mathbf{r}, \mathbf{v}_t)$ . See Section H.3 in the appendix for proof.

**COROLLARY 2.** *Let  $n$  be the number of bidders. Assume that the maximum valuation is normalized to one. Then, for the problem of maximizing personalized reserve prices in the online full information setting, there exists a learning algorithm that obtains  $O\left(nT^{1/2} \log^{1/2} T\right)$   $\frac{1}{2}$ -regret, where  $T$  is the number of auctions. Furthermore, in the online bandit setting, there exists a learning algorithm that obtains  $O\left(n^{3/5}T^{4/5} \log^{1/3}(nT)\right)$   $\frac{1}{2}$ -regret. Here, the benchmark in the regret bounds is  $\frac{1}{2} \max_{\mathbf{r} \in [0, 1]^n} \sum_{t=1}^T f(\mathbf{r}, \mathbf{v}_t)$ .*

## Acknowledgments

N.G. was supported in part by the Young Investigator Program (YIP) Award from the Office of Naval Research (ONR) N00014-21-1-2776 and the MIT Research Support Award. We thank Tim Roughgarden, Dimitris Bertsimas, and Amin Karbasi for their insightful comments during this work.

## References

- Jacob Abernethy, Elad E Hazan, and Alexander Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. In *21st Annual Conference on Learning Theory, COLT 2008*, pages 263–273, 2008.
- Jacob Abernethy, Peter L Bartlett, and Elad Hazan. Blackwell approachability and no-regret learning are equivalent. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 27–46, 2011.
- Shabbir Ahmed and Alper Atamtürk. Maximizing a class of submodular utility functions. *Mathematical programming*, 128(1-2):149–169, 2011.
- Saeed Alaei, Jason Hartline, Rad Niazadeh, Emmanouil Pountourakis, and Yang Yuan. Optimal auctions vs. anonymous pricing. *Games and Economic Behavior*, 118:494–510, 2019.
- Ali Aouad and Danny Segev. Display optimization for vertically differentiated locations under multinomial logit preferences. *Management Science*, 67(6):3519–3550, 2021.

- Arash Asadpour, Rad Niazadeh, Amin Saberi, and Ali Shameli. Ranking an assortment of products via sequential submodular optimization. *Available at SSRN*, 2020.
- Susan Athey and Glenn Ellison. Position auctions with consumer search. *The Quarterly Journal of Economics*, 126(3):1213–1270, 2011.
- Jean-Yves Audibert, Sébastien Bubeck, and Gábor Lugosi. Regret in online combinatorial optimization. *Mathematics of Operations Research*, 39(1):31–45, 2014.
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- Francis Bach. Submodular functions: from discrete to continuous domains. *Mathematical Programming*, 175(1):419–459, 2019.
- Hedyeh Beyhaghi, Negin Golrezaei, Renato Paes Leme, Martin Pal, and Balasubramanian Sivan. Improved approximations for free-order prophets and second-price auctions. *arXiv preprint arXiv:1807.03435*, 2018.
- Andrew An Bian, Baharan Mirzasoleiman, Joachim Buhmann, and Andreas Krause. Guaranteed non-convex optimization: Submodular maximization over continuous domains. In *Artificial Intelligence and Statistics*, pages 111–120. PMLR, 2017.
- David Blackwell. An analog of the minimax theorem for vector payoffs. *Pacific Journal of Mathematics*, 6(1):1–8, 1956.
- Sébastien Bubeck, Nicolo Cesa-Bianchi, and Sham M Kakade. Towards minimax policies for online linear optimization with bandit feedback. In *Conference on Learning Theory*, pages 41–1, 2012.
- Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- Niv Buchbinder and Moran Feldman. Deterministic algorithms for submodular maximization problems. *ACM Transactions on Algorithms (TALG)*, 14(3):32, 2018.
- Niv Buchbinder, Moran Feldman, Joseph Seffi, and Roy Schwartz. A tight linear time (1/2)-approximation for unconstrained submodular maximization. *SIAM Journal on Computing*, 44(5):1384–1402, 2015.
- Gruia Calinescu, Chandra Chekuri, Martin Pal, and Jan Vondrák. Maximizing a monotone submodular function subject to a matroid constraint. *SIAM Journal on Computing*, 40(6):1740–1766, 2011.
- Nicolo Cesa-Bianchi and Gabor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.
- Nicolo Cesa-Bianchi and Gábor Lugosi. Combinatorial bandits. *Journal of Computer and System Sciences*, 78(5):1404–1422, 2012.
- Nicolo Cesa-Bianchi, Claudio Gentile, and Yishay Mansour. Regret minimization for reserve prices in second-price auctions. *IEEE Transactions on Information Theory*, 61(1):549–564, 2014.

- Lin Chen, Christopher Harshaw, Hamed Hassani, and Amin Karbasi. Projection-free online optimization with stochastic gradient: From convexity to submodularity. In *International Conference on Machine Learning*, pages 814–823. PMLR, 2018a.
- Lin Chen, Hamed Hassani, and Amin Karbasi. Online continuous submodular maximization. In *International Conference on Artificial Intelligence and Statistics*, pages 1896–1905. PMLR, 2018b.
- Lin Chen, Mingrui Zhang, Hamed Hassani, and Amin Karbasi. Black box submodular maximization: Discrete and continuous settings. In *International Conference on Artificial Intelligence and Statistics*, pages 1058–1070. PMLR, 2020.
- Wei Chen, Yajun Wang, and Yang Yuan. Combinatorial multi-armed bandit: General framework and applications. In *International Conference on Machine Learning*, pages 151–159, 2013.
- Richard Combes, Mohammad Sadegh Talebi Mazraeh Shahi, Alexandre Proutiere, et al. Combinatorial bandits revisited. In *Advances in Neural Information Processing Systems*, pages 2116–2124, 2015.
- Mahsa Derakhshan, Negin Golrezaei, and Renato Paes Leme. Lp-based approximation for personalized reserve prices. In *Proceedings of the 2019 ACM Conference on Economics and Computation*, pages 589–589, 2019.
- Mahsa Derakhshan, Negin Golrezaei, Vahideh Manshadi, and Vahab Mirrokni. Product ranking on online platforms. In *Proceedings of the 21st ACM Conference on Economics and Computation*, pages 459–459, 2020.
- Shahar Dobzinski and Michael Schapira. An improved approximation algorithm for combinatorial auctions with submodular bidders. In *Proceedings of the seventeenth annual ACM-SIAM symposium on Discrete algorithm*, pages 1064–1073, 2006.
- Miroslav Dudík, Nika Haghtalab, Haipeng Luo, Robert E Schapire, Vasilis Syrgkanis, and Jennifer Wortman Vaughan. Oracle-efficient online learning and auction design. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 528–539. IEEE, 2017.
- Eyal Even-Dar, Robert Kleinberg, Shie Mannor, and Yishay Mansour. Online learning for global cost functions. In *COLT*, 2009.
- Uriel Feige, Vahab S Mirrokni, and Jan Vondrák. Maximizing non-monotone submodular functions. *SIAM Journal on Computing*, 40(4):1133–1153, 2011.
- Moran Feldman, Joseph Naor, and Roy Schwartz. A unified continuous greedy algorithm for submodular maximization. In *2011 IEEE 52nd Annual Symposium on Foundations of Computer Science*, pages 570–579. IEEE, 2011.
- Kris J Ferreira, Sunanda Parthasarathy, and Shreyas Sekar. Learning to rank an assortment of products. *Management Science*, 2021.
- Abraham D Flaxman, Adam Tauman Kalai, and H Brendan McMahan. Online convex optimization in the bandit setting: gradient descent without a gradient. In *Proceedings of the sixteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 385–394, 2005.

- Dan Garber. Efficient online linear optimization with approximation algorithms. *Mathematics of Operations Research*, 46(1):204–220, 2021.
- Negin Golrezaei, Adel Javanmard, and Vahab Mirrokni. Dynamic incentive-aware learning: Robust pricing in contextual auctions. *Operations Research*, 69(1):297–314, 2021a.
- Negin Golrezaei, Vahideh Manshadi, Jon Schneider, and Shreyas Sekar. Learning product rankings robust to fake users. In *Proceedings of the 22nd ACM Conference on Economics and Computation*, pages 560–561, 2021b.
- Jason D. Hartline and Tim Roughgarden. Simple versus optimal mechanisms. In *Proceedings 10th ACM Conference on Electronic Commerce (EC-2009), Stanford, California, USA, July 6–10, 2009*, pages 225–234, 2009.
- Hamed Hassani, Mahdi Soltanolkotabi, and Amin Karbasi. Gradient methods for submodular maximization. In *Advances in Neural Information Processing Systems*, pages 5841–5851, 2017.
- Hamed Hassani, Amin Karbasi, Aryan Mokhtari, and Zebang Shen. Stochastic conditional gradient++:(non) convex minimization and continuous submodular maximization. *SIAM Journal on Optimization*, 30(4):3315–3344, 2020.
- Elad Hazan and Zohar Karnin. Volumetric spanners: an efficient exploration basis for learning. *The Journal of Machine Learning Research*, 17(1):4062–4095, 2016.
- Elad Hazan and Tomer Koren. The computational power of optimization in online learning. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 128–141. ACM, 2016.
- Elad Hazan, Wei Hu, Yuanzhi Li, and Zhiyuan Li. Online improper learning with an approximation oracle. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 5657–5665, 2018.
- Sham M Kakade, Adam Tauman Kalai, and Katrina Ligett. Playing games with approximation algorithms. *SIAM Journal on Computing*, 39(3):1088–1106, 2009.
- Adam Kalai and Santosh Vempala. Efficient algorithms for online decision problems. *Journal of Computer and System Sciences*, 71(3):291–307, 2005.
- David Kempe, Jon Kleinberg, and Éva Tardos. Maximizing the spread of influence through a social network. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 137–146, 2003.
- A Gürhan Kök, Marshall L Fisher, and Ramnath Vaidyanathan. Assortment planning: Review of literature and industry practice. In *Retail supply chain management*, pages 99–153. Springer, 2008.
- Joon Kwon and Vianney Perchet. Online learning and blackwell approachability with partial monitoring: optimal convergence rates. In *Artificial Intelligence and Statistics*, pages 604–613, 2017.
- Ehud Lehrer. Approachability in infinite dimensional spaces. *International Journal of Game Theory*, 31(2): 253–268, 2003.

- Shie Mannor and Nahum Shimkin. Online learning with variable stage duration. In *International Conference on Computational Learning Theory*, pages 408–422. Springer, 2006.
- Shie Mannor, Vianney Perchet, and Gilles Stoltz. Robust approachability and regret minimization in games with partial monitoring. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 515–536, 2011.
- Shie Mannor, Vianney Perchet, and Gilles Stoltz. Set-valued approachability and online learning with partial monitoring. *The Journal of Machine Learning Research*, 15(1):3247–3295, 2014.
- Emanuel Milman. Approachable sets of vector payoffs in stochastic games. *Games and Economic Behavior*, 56(1):135–147, 2006.
- Marko Mitrovic, Moran Feldman, Andreas Krause, and Amin Karbasi. Submodularity on hypergraphs: From sets to sequences. In *International Conference on Artificial Intelligence and Statistics*, pages 1177–1184. PMLR, 2018.
- Aryan Mokhtari, Hamed Hassani, and Amin Karbasi. Stochastic conditional gradient methods: From convex minimization to submodular maximization. *Journal of machine learning research*, 2020.
- George L Nemhauser, Laurence A Wolsey, and Marshall L Fisher. An analysis of approximations for maximizing submodular set functions—i. *Mathematical programming*, 14(1):265–294, 1978.
- Rad Niazadeh, Tim Roughgarden, and Joshua Wang. Optimal algorithms for continuous non-monotone submodular and dr-submodular maximization. In *Advances in Neural Information Processing Systems*, pages 9594–9604, 2018.
- Tim Roughgarden and Joshua R Wang. An optimal learning algorithm for online unconstrained submodular maximization. In *Conference On Learning Theory*, pages 1307–1325, 2018.
- Tim Roughgarden and Joshua R Wang. Minimizing regret with multiple reserves. *ACM Transactions on Economics and Computation (TEAC)*, 7(3):1–18, 2019.
- Shai Shalev-Shwartz. Online learning: Theory, algorithms, and applications. *PhD Thesis*, 2007.
- Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194, 2012.
- Xavier Spinat. A necessary and sufficient condition for approachability. *Mathematics of Operations Research*, 27(1):31–44, 2002.
- M Streeter and D Golovin. An online algorithm for maximizing submodular functions (technical report cmu-cs-07-171), 2007.
- Matthew Streeter and Daniel Golovin. An online algorithm for maximizing submodular functions. In *Advances in Neural Information Processing Systems*, pages 1577–1584, 2008.
- Nguyen Kim Thang and Abhinav Srivastav. Online non-monotone dr-submodular maximization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 9868–9876, 2021.

- Sebastian Tschiatschek, Adish Singla, and Andreas Krause. Selecting sequences of items via submodular maximization. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- Taishi Uchiya, Atsuyoshi Nakamura, and Mineichi Kudo. Algorithms for adversarial bandit problems with multiple plays. In *International Conference on Algorithmic Learning Theory*, pages 375–389. Springer, 2010.
- Raluca Mihaela Ursu. The power of rankings: Quantifying the effect of rankings on online consumer search and purchase decisions. *Browser Download This Paper*, 2016.
- Nicolas Vieille. Weak approachability. *Mathematics of operations research*, 17(4):781–791, 1992.
- Jan Vondrák. Optimal approximation for the submodular welfare problem in the value oracle model. In *Proceedings of the fortieth annual ACM symposium on Theory of computing*, pages 67–74, 2008.
- H Martin Weingartner. *Mathematical programming and the analysis of capital budgeting problems*. Markham Publishing Company, 1967.
- Laurence A Wolsey. An analysis of the greedy algorithm for the submodular set covering problem. *Combinatorica*, 2(4):385–393, 1982.
- Andrew Chi-Chin Yao. Probabilistic computations: Toward a unified measure of complexity. In *18th Annual Symposium on Foundations of Computer Science (sfcs 1977)*, pages 222–227. IEEE, 1977.
- Mingrui Zhang, Lin Chen, Hamed Hassani, and Amin Karbasi. Online continuous submodular maximization: From full-information to bandit feedback. In *Advances in Neural Information Processing Systems*, pages 9206–9217, 2019.
- Mingrui Zhang, Zebang Shen, Aryan Mokhtari, Hamed Hassani, and Amin Karbasi. One sample stochastic frank-wolfe. In *International Conference on Artificial Intelligence and Statistics*, pages 4012–4023. PMLR, 2020.
- Julian Zimmert, Haipeng Luo, and Chen-Yu Wei. Beating stochastic and adversarial semi-bandits optimally and simultaneously. In *International Conference on Machine Learning*, pages 7683–7692. PMLR, 2019.

## Appendix A: Numerical Studies

In this section, we evaluate our two algorithms ONLINE-IG (Algorithm 2) and BANDIT-IG (Algorithm 4) for two of the applications we presented in Section 6: (i) product ranking problem and (ii) maximizing multiple reserves in second price auction. Recall that ONLINE-IG is defined/designed for the full information setting while BANDIT-IG is defined/designed for the bandit setting.

### A.1. Numerical Studies for Product Ranking Problem

**Simulation setting.** In our setting, we have  $n = 10$  items where  $\mathcal{N} = \{1, 2, \dots, 10\}$  is the universe of items,  $T = 20,000$  customers (rounds), and  $k = 2$  types appearing equally likely. We consider multinomial logit (MNL) choice models for customers. Given this choice modeling, a type  $u$  is characterized by a patience level  $\theta_u$  and a set of weights over items  $\mathbf{w}_u \in \mathbb{R}^n$  such that given a ranked list of products  $\boldsymbol{\pi} = ([\boldsymbol{\pi}]_1, [\boldsymbol{\pi}]_2, \dots, [\boldsymbol{\pi}]_n)$ , they inspect the first  $\theta_u$  products and clicks with probability  $\kappa_u(\{[\boldsymbol{\pi}]_1, [\boldsymbol{\pi}]_2, \dots, [\boldsymbol{\pi}]_{\theta_u}\}) = \frac{w_{u, [\boldsymbol{\pi}]_1} + w_{u, [\boldsymbol{\pi}]_2} + \dots + w_{u, [\boldsymbol{\pi}]_{\theta_u}}}{1 + w_{u, [\boldsymbol{\pi}]_1} + w_{u, [\boldsymbol{\pi}]_2} + \dots + w_{u, [\boldsymbol{\pi}]_{\theta_u}}}$ , following the MNL model. Let  $f(\boldsymbol{\pi})$  denote the expected market share when ranking  $\boldsymbol{\pi}$  is offered.<sup>19</sup>

We change the parameters  $(\theta_{u,t}, \mathbf{w}_{u,t})_{u=1,2}$  every 500 rounds (one episode), and keep them constant in the same episode. At the end of an even episode, we set patience levels  $\theta_{1,t} = \theta_{2,t} = 10$ , while at the end of an odd episode, we set both of them to be 5. Furthermore, at the end of an even episode, we choose 3 items uniformly at random from  $\{1, 2, \dots, 5\}$  for both types, and put equal weight of  $1/3$ , on each of the chosen items, and zero weights on the remaining items. On the other hand, at the end of an odd episode, we choose 3 items uniformly at random from  $\{6, \dots, 10\}$ , then put equal weights on the chosen items, and 0 for the remaining items. This setting allows us to model the adversarial nature of the environment.

**Implementation of ONLINE-IG algorithm.** We replace AlgB with the Multiplicative Weights or Hedge algorithm with the learning rate of  $\epsilon_t = \sqrt{\frac{1}{t}}$  with  $n = 10$  arms (representing items). The Hedge algorithm is a no-regret adversarial learning algorithm that keeps a weight over the different arms at each time step, and updates those according to the observed feedback. Recall that in the product ranking problem, subproblem  $i$  corresponds to determining the item for the  $i^{\text{th}}$  position in the ranking. Then, the rewards we give to the Hedge algorithm that corresponds to subproblem  $i$  is the marginal market share of adding item  $j$  to the top  $(i - 1)$  items, for all  $j \in \mathcal{N}$ . Specifically, suppose that in round  $t$ , the top  $(i - 1)$  items are  $([\boldsymbol{\pi}]_1, [\boldsymbol{\pi}]_2, \dots, [\boldsymbol{\pi}]_{i-1})$ . Then, the set of rewards we feed to hedge algorithm in our setting is  $f(\{[\boldsymbol{\pi}]_1, [\boldsymbol{\pi}]_2, \dots, [\boldsymbol{\pi}]_{i-1}, j\}) - f(\{[\boldsymbol{\pi}]_1, [\boldsymbol{\pi}]_2, \dots, [\boldsymbol{\pi}]_{i-1}\})$  for  $j = 1, 2, \dots, 10$ .

**Implementation of BANDIT-IG algorithm.** For BANDIT-IG, we modify the Hedge algorithm to implement AlgBB, and we call this the Hedge-bandit algorithm. The Hedge-bandit algorithm has exploration and exploitation rounds and only updates its weights in the exploration rounds. Following AlgBB, the Hedge-bandit algorithm explores with probability  $q = (\frac{n}{T})^{1/3}$ ; see Algorithm 3 for details. In each round  $t \in [T]$  and for each subproblem  $i \in [10]$ , the algorithm decides to explore with probability  $q$ , and continues to the next subproblem with probability  $(1 - q)$ . Suppose that in round  $t$ , the first  $(i - 1)$  subproblems do not enter the exploration round, and the ranking so far is  $\boldsymbol{\pi}^{(i-1)} = ([\boldsymbol{\pi}]_1, [\boldsymbol{\pi}]_2, \dots, [\boldsymbol{\pi}]_{i-1})$ . For position  $i$ , in the

<sup>19</sup> Here, we have  $f(\boldsymbol{\pi}) = \frac{1}{2} \sum_{u=1}^2 \kappa_u(\{[\boldsymbol{\pi}]_1, [\boldsymbol{\pi}]_2, \dots, [\boldsymbol{\pi}]_{\theta_u}\})$  as both types are equally likely.

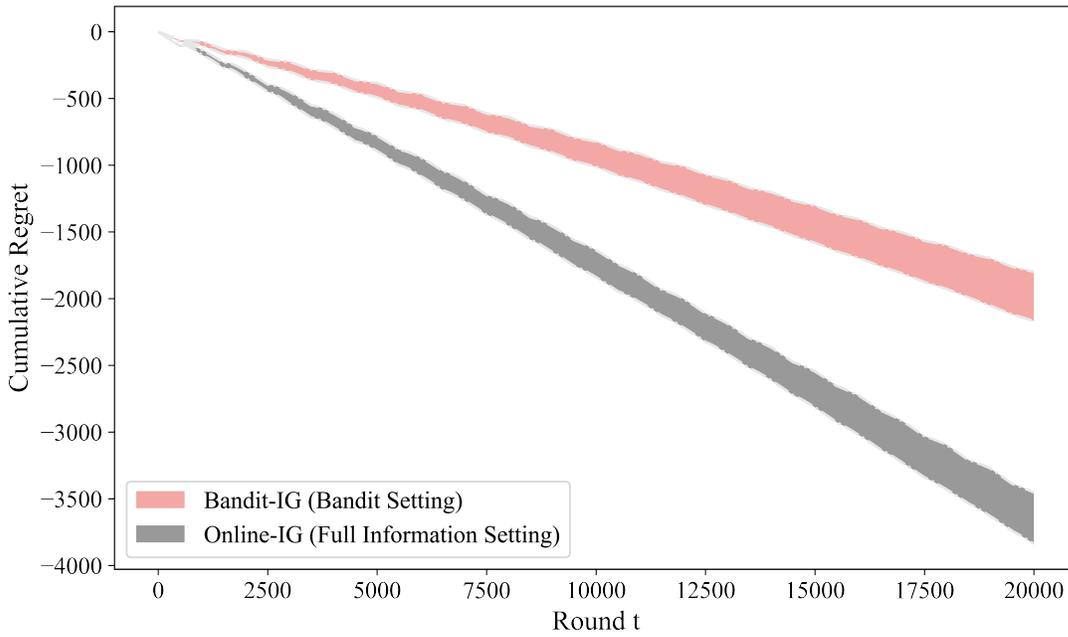
exploration round, we choose an item  $j \in \mathcal{N}$  uniformly at random, then with probability 0.5, the algorithm outputs the ranking  $\boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_i = ([\boldsymbol{\pi}]_1, [\boldsymbol{\pi}]_2, \dots, [\boldsymbol{\pi}]_{i-1}, j)$ , observes its market share  $f(\boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_i)$ , and feeds  $2n/ql f(\boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_i)$  back to the Hedge-bandit algorithm. Furthermore, with probability 0.5, it outputs the ranking  $\boldsymbol{\pi}^{(i-1)} = ([\boldsymbol{\pi}]_1, [\boldsymbol{\pi}]_2, \dots, [\boldsymbol{\pi}]_{i-1})$ , observes its market share  $f(\boldsymbol{\pi}^{(i-1)})$ , and feeds  $-2n/ql f(\boldsymbol{\pi}^{(i-1)})$  to the Hedge-bandit algorithm. If we consider an  $n$ -dimensional vector with the resulting randomized feedback in its  $j^{\text{th}}$  coordinate and zero elsewhere, this vector is an unbiased estimate of the marginal market share vector  $[f(\boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_i) - f(\boldsymbol{\pi}^{(i-1)})]_{j \in [n]}$ . This is based on the feedback stated in Equation (39),  $n(\theta_j \mathbf{1}_n - \mathbf{e}_j)f(\boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_i)$ , which also incorporates  $\boldsymbol{\theta}$ , the current weight over possible items for subproblem  $i$ .

**Regret.** Figure 1a shows the cumulative regret of ONLINE-IG and BANDIT-IG algorithms over time. In computing the regret of our algorithms, we use 1/2 of the optimal market share as the regret benchmark. Note that the optimal market share is the highest possible market share among all rankings assuming that we know all the parameters (weights and patience windows for all types) beforehand.<sup>20</sup> We choose the factor 1/2 because the greedy algorithm that we based our online algorithm on is a 1/2-approximation of the optimal value. In Figure 1a, we consider 50 instances where each problem has a unique set of parameters (i.e.  $(\theta_{u,t}, \mathbf{w}_{u,t})_{u \in [2], t \in [T]}$ ). For each instance, we take the average performance of both ONLINE-IG and BANDIT-IG over 10 runs. Then, we calculate the regret of these algorithms using 1/2 of the optimal values as the benchmark. As we observe from Figure 1a, BANDIT-IG has a bigger regret than that of ONLINE-IG. In addition, the regret of both algorithms are negative, implying that they do better than half of the optimal market share.

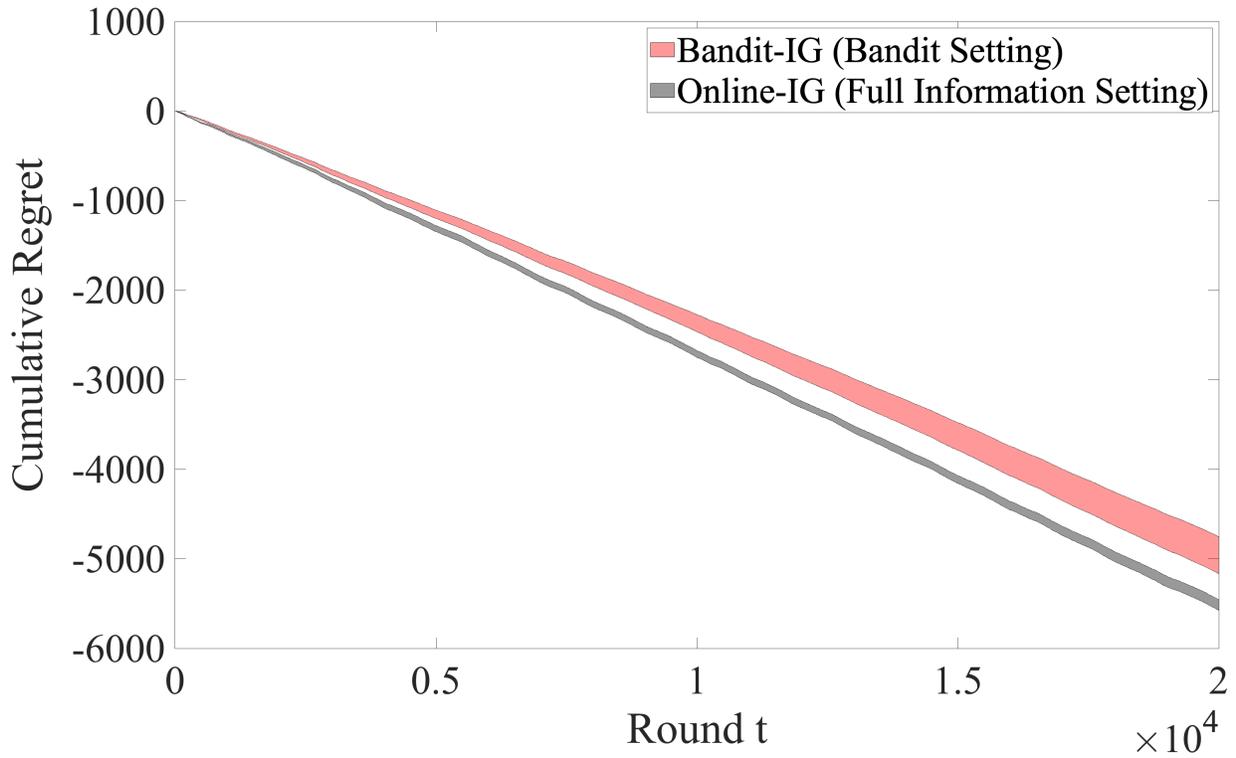
## A.2. Numerical Studies for Maximizing Multiple Reserves in Second Price Auctions

**Simulation setting.** We consider a setting with  $n = 3$  buyers and  $T = 20,000$  auctions. In each round  $t$ , to generate the value of each buyer  $i \in [3]$ , denoted by  $v_{i,t}$ , we first generate an independent random variable  $\hat{v}_{i,t}$  form a log-normal distribution with parameters  $(\mu_{i,t}, \sigma_{i,t})$ , where  $\mu_{i,t} \in [0, 0.5]$  and  $\sigma_{i,t} \in [0, 0.2]$ . We then set  $v_{i,t} = \hat{v}_{i,t} \cdot x_{i,t}$ , where  $x_{i,t} \in \{0, 1\}$ . Here,  $x_{i,t}$ , which is independently drawn from a Bernoulli distribution with success probability of  $p_{i,t}$ , models the fact that not all the buyers participate in all the auctions. In our setting, parameters  $(\mu_{i,t}, \sigma_{i,t}, p_{i,t})_{i \in [3]}$  remain fixed during an episode with the length of 500 rounds, but keep changing across episodes, modeling the adversarial nature of the environment. In addition, within odd episodes, we set  $\mu_{i,t}$  and  $\sigma_{i,t}$  to  $\mu_{i,1}$  and  $\sigma_{i,1}$ , and similarly, within even episodes, we set  $\mu_{i,t}$  and  $\sigma_{i,t}$  to  $\mu_{i,2}$  and  $\sigma_{i,2}$ . Here, we independently draw  $\mu_{i,j}$  and  $\sigma_{i,j}$ ,  $i \in [n]$  and  $j \in [2]$ , from a uniform distribution in the range of  $[0, 0.5]$  and  $[0, 0.2]$ , respectively. Similarly, at the end of each episode, to set  $p_{i,t}$ ,  $i \in [n]$ , we draw a uniform random variable in the range of  $[0, 1]$ . (Note that  $p_{i,t}$  does not follow the alternating episodic pattern that  $\mu_{i,t}$  and  $\sigma_{i,t}$  follow.) Finally, for reserve prices, we assume that reserve prices belongs to set  $\mathcal{R} = \{\rho_1, \dots, \rho_{20}\}$ , where  $\rho_1 = 0.1$ ,  $\rho_{20}$  is the maximum buyers' value across all the auctions, and  $\frac{\rho_i}{\rho_{i+1}}$ ,  $i \in [19]$ , is equal to some constant  $k$ .

<sup>20</sup> The optimal ranking is obtained by enumerating all possible rankings.



(a) Product Ranking



(b) Reserve Price Optimization

Figure 1 Average cumulative regret of Online-IG and Bandit-IG over time in the product ranking and reserve price optimization problems. The width of the curves is equal to two times the standard error of the regret bounds across 50 problem instances.

**Implementation of ONLINE-IG and BANDIT-IG algorithms.** To implement ONLINE-IG, we again use the Hedge algorithm with the learning rate of  $\epsilon_t = \sqrt{\frac{10}{t}}$  and the number of arms of 20. (Recall the number of feasible reserve prices is equal to 20.) At the end of each round, for every subproblem  $i \in [3]$ , we pass  $q^{(i)}(r)$ ,  $r \in \mathcal{R}$  to the Hedge algorithm as the reward of reserve price  $r$ , where  $\mathcal{R}$  is the set of feasible reserve prices. (See the definition of  $q^{(i)}(r)$  in Algorithm 6.) For BANDIT-IG, we again rely on the Hedge algorithm to implement AlgBB. The Hedge-bandit algorithm explores with probability  $q = (2\frac{n}{T})^{1/3}$ . We now explain how the Hedge-bandit algorithm performs during an exploration round. Suppose that in some round  $t$ , we would like to decide about the reserve price of some buyer  $i$  and the first  $i - 1$  buyers (subproblems) did not enter an exploration round. If buyer/subproblem  $i$  enters an exploration round, given the sampling device presented in Section H.2 (see Equation (40)), we first choose one of the reserve prices in  $\mathcal{R}$  uniformly at random. Let us denote this reserve price by  $\hat{\rho}$ . Then, with probability  $1/2$ , we set the reserve price of all the buyers (including buyer  $i$ ) equal to  $\hat{\rho}$ . Otherwise, we set the reserve price of all the buyers, except buyer  $i$ , equal to  $\hat{\rho}$ . The reserve price of buyer  $i$  is then set to 0. In the former case, we pass  $2mf(\hat{\rho}\mathbf{1}_n, \mathbf{v}_t)$  to the Hedge-bandit algorithm as the reward of arm  $\hat{\rho}$ . The reward of other arms is set to zero. Here,  $m = 20$  is the number of feasible reserve prices and  $f(\hat{\rho}\mathbf{1}_n, \mathbf{v}_t)$  is the revenue of the second price auction with reserve prices  $\hat{\rho}\mathbf{1}_n$  when buyers' values are  $\mathbf{v}_t$ . In the latter case, we pass  $-2mf(\hat{\rho}(\mathbf{1}_n - \mathbf{e}_i), \mathbf{v}_t)$  to the Hedge-bandit algorithm as the reward of arm  $\hat{\rho}$ . The reward of other arms is set to zero.<sup>21</sup>

**Regret.** Figure 1b shows the cumulative regret of ONLINE-IG and BANDIT-IG over time. We use  $1/2$  of the optimal revenue as our benchmark in our regret calculation, where we compute the optimal revenue by enumerating over all possible vector of reserve prices. (Recall that our online algorithms are transformation of the  $1/2$ -approximate greedy algorithm of Roughgarden and Wang (2019).) In Figure 1b, we consider 50 problem instances where each problem instance has a unique set of parameters (i.e.,  $(\mu_{i,t}, \sigma_{i,t}, p_{i,t})_{i \in [n], t \in [T]}$ ), as well as, auction bids/values generated from these parameters. For each problem instance, we run ONLINE-IG and BANDIT-IG 10 times and compute their average performance. We observe that ONLINE-IG outperforms BANDIT-IG, as expected. Further, similar to our results for the product ranking problem, the regret of both algorithms is negative.

## Appendix B: Proofs and Remarks of Section 2.3

### B.1. Equivalent criteria for approachability

Interestingly, there are other structural conditions that are equivalent to approachability. For example, the original proof of the Blackwell approachability theorem (Blackwell 1956) uses a condition called “halfspace-satisfiability”. The following proposition summarizes all the known equivalences.

<sup>21</sup> Considering the fact that all the subproblems need the value of  $f(r\mathbf{1}_n, \mathbf{v}_t)$ , for some  $r \in \mathcal{R}$ , during their exploration rounds, in our implementation, we allow different subproblems to learn from each other. In particular, when subproblem  $i$  explores, with probability  $1/(n+1)$ , we choose reserve price of  $\hat{\rho}\mathbf{1}_n$  and with probability  $n/(n+1)$ , we choose reserve price of  $\hat{\rho}(\mathbf{1}_n - \mathbf{e}_i)$ , where  $\hat{\rho}$  is randomly chosen reserve price by subproblem  $i$ . Then, when  $\hat{\rho}\mathbf{1}_n$  is chosen, we update the weights of all the Hedge-bandit algorithms, not only those of subproblem  $i$ . To do so, we pass the reward of  $(n+1) \cdot m \cdot f(\hat{\rho}\mathbf{1}_n, \mathbf{v}_t)$  to the Hedge-bandit algorithms. When reserve price of  $\hat{\rho}(\mathbf{1}_n - \mathbf{e}_i)$  is chosen, we only update the Hedge-bandit algorithm of subproblem  $i$  by passing the reward of  $-\frac{n+1}{n} \cdot m \cdot f(\hat{\rho}(\mathbf{1}_n - \mathbf{e}_i), \mathbf{v}_t)$ .

PROPOSITION 1 (**Satisfiable/Halfspace-Satisfiable/Response-Satisfiable (Abernethy et al. 2011)**).

The following conditions are all equivalent to the approachability condition (Definition 3):

1. A target set  $S$  is satisfiable in the Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  if there exists a player 1's action  $\mathbf{x} \in \mathcal{X}$  such that for every player 2's action  $\mathbf{y} \in \mathcal{Y}$ , the vector payoff falls into the target set, that is  $\mathbf{p}(\mathbf{x}, \mathbf{y}) \in S$ .
2. A target set  $S$  is halfspace-satisfiable in the Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  if for every halfspace  $H \supseteq S$ ,  $H$  is satisfiable.
3. A target set  $S$  is response-satisfiable in the Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  if for every player 2's action  $\mathbf{y} \in \mathcal{Y}$ , there exists a player 1's action  $\mathbf{x} \in \mathcal{X}$  such that the vector payoff falls into the target set, that is  $\mathbf{p}(\mathbf{x}, \mathbf{y}) \in S$ .

## B.2. Proof of Theorem 1

*Proof.* The proof for the only if direction relies on the fact that the  $\ell_\infty$ -distance between the average payoff and  $S$  is vanishing as  $T \rightarrow +\infty$  since  $S$  is  $o(1)$ -approachable. Suppose that  $S$  is not response satisfiable, then there exists player 2's action  $\mathbf{y}_0 \in \mathcal{Y}$  such that for every player 1's action  $\mathbf{x} \in \mathcal{X}$ , the payoff  $\mathbf{p}(\mathbf{x}, \mathbf{y}_0)$  is not in  $S$ . Consider the set  $U := \{\mathbf{p}(\mathbf{x}, \mathbf{y}_0) : \mathbf{x} \in \mathcal{X}\}$ . Because the payoff  $\mathbf{p}$  is biaffine and  $\mathcal{X}$  is convex and compact, so is  $U$ , hence  $\inf_{\mathbf{u} \in U} d_\infty(\mathbf{u}, S) = d_\infty(\mathbf{p}(\underline{\mathbf{x}}, \mathbf{y}_0), S)$  for some  $\underline{\mathbf{x}} \in \mathcal{X}$ . As  $\mathbf{p}(\underline{\mathbf{x}}, \mathbf{y}_0) \notin S$ ,  $\beta = d_\infty(\mathbf{p}(\underline{\mathbf{x}}, \mathbf{y}_0), S) > 0$ . When player 2 always plays  $\mathbf{y}_0$ , we know that the  $\ell_\infty$  distance between the average payoff and  $S$  should converge to zero as  $S$  is  $o(1)$ -approachable. At the same time, it is at least  $\beta$ , a contradiction.

To prove the if direction, we first show a reduction from Blackwell approachability to Online Linear Optimization (OLO) by showing that we can upper bound the  $\ell_\infty$  distance between the average payoff and the target set in a Blackwell approachability problem with the regret of the corresponding OLO instance. Then, we bound the regret of the OLO problem from above in terms of the  $\ell_\infty$  norm of the payoff  $D(\mathbf{p})$  (because of our desired bound), the number of rounds  $T$ , and the dimension of the payoff function  $d$ . We assume that  $S$  is a cone throughout the proof, which is not an issue because we can always lift a convex set to a cone in one dimension higher while not perturbing the distances by more than a factor of 2.

*Blackwell approachability reduces to OLO.* In an OLO problem, a player is given a compact convex decision set  $\mathcal{K} \subset \mathbb{R}^d$ , and have to decide on a sequence of actions  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_T \in \mathcal{K}$ . In round  $t$ , after the player decides on an action  $\mathbf{w}_t$ , Nature reveals a loss vector  $\mathbf{l}_t$  and the player pays  $\langle \mathbf{l}_t, \mathbf{w}_t \rangle$ . The player observes the loss vector  $\mathbf{l}_t$  in each round (full-information setting) and aims to minimize his cost. We want to construct a learning algorithm  $\mathcal{L}$ , such that, for any sequence of loss vectors  $\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_T \in \mathbb{R}^d$ , the algorithm outputs  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_T \in \mathcal{K}$  that attains a small regret, i.e.  $\sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w}_t \rangle - \min_{\mathbf{w} \in \mathcal{K}} \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w} \rangle \leq o(T)$ . Abernethy et al. (2011) show that we can efficiently obtain an algorithm for a Blackwell approachability problem from an algorithm for its corresponding OLO problem. Specifically, we have the following lemma.

LEMMA 1. (*Abernethy et al. (2011)*) Given a Blackwell instance  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$ , and a cone  $S$  such that  $S$  is response-satisfiable, we can construct an OLO problem with  $\mathcal{K} = S^\circ \cap B_2(1)$ <sup>22</sup> and  $\mathbf{l}_t = -\mathbf{p}(\mathbf{x}_t, \mathbf{y}_t)$  for all  $t$ , such that, if the OLO learning algorithm returns  $\mathbf{w}_t$  in round  $t$ , we can convert it into  $\mathbf{x}_t \in \mathcal{X}$  where

$$d_2\left(\frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S\right) \leq \frac{1}{T} \left( \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w}_t \rangle - \min_{\mathbf{w} \in \mathcal{K}} \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w} \rangle \right).$$

*Proof of Lemma 1.* This lemma was proved in [Abernethy et al. \(2011\)](#), but we include the proof here for completion. Notice that, for any  $\mathbf{x} \in \mathbb{R}^d$  and convex cone  $S \subseteq \mathbb{R}^d$ , the distance from  $\mathbf{x}$  to  $S$  can be written as

$$d_\infty(\mathbf{x}, S) = \max_{\mathbf{w} \in S^\circ, \|\mathbf{w}\| \leq 1} \langle \mathbf{w}, \mathbf{x} \rangle \quad (4)$$

because

$$d_\infty(\mathbf{x}, S) = \|\mathbf{x} - \pi_S(\mathbf{x})\|_2 \geq \|\mathbf{w}\| \|\mathbf{x} - \pi_S(\mathbf{x})\| \geq \langle \mathbf{w}, \mathbf{x} - \pi_S(\mathbf{x}) \rangle \geq \langle \mathbf{w}, \mathbf{x} \rangle,$$

where  $\pi_S(\mathbf{x})$  denotes the projection of  $\mathbf{x}$  onto  $S$ , and when  $\mathbf{w} = \frac{\mathbf{x} - \pi_S(\mathbf{x})}{\|\mathbf{x} - \pi_S(\mathbf{x})\|_2}$ , we have equality, i.e.  $\langle \mathbf{w}, \mathbf{x} \rangle = d_2(\mathbf{x}, S)$ . To construct a mapping from the output of the OLO algorithm  $\mathbf{w}_t$  to  $\mathbf{x}_t$  for the Blackwell game, we utilize the halfspace oracle for the Blackwell problem (see [Proposition 1](#)). Specifically, we pick  $\mathbf{x}_t$  such that  $\mathbf{p}(\mathbf{x}_t, \mathbf{y}_t) \in H_{\mathbf{w}_t}$  for all  $\mathbf{y} \in \mathcal{Y}$ , where  $H_{\mathbf{w}_t} = \{\mathbf{x} : \langle \mathbf{w}_t, \mathbf{x} \rangle \leq 0\}$  is a halfspace that contains  $S$  ( $H_{\mathbf{w}_t}$  contains  $S$  because its normal,  $\mathbf{w}_t$ , is in  $S^\circ$ ). This gives us the following guarantee

$$\begin{aligned} d_2\left(\frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S\right) &\stackrel{(1)}{=} \max_{\mathbf{w} \in \mathcal{K}} \left\langle \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), \mathbf{w} \right\rangle = \frac{1}{T} \max_{\mathbf{w} \in \mathcal{K}} \left( - \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w} \rangle \right) \\ &\stackrel{(2)}{\leq} \frac{1}{T} \left( \sum_{t=1}^T \langle -\mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), \mathbf{w}_t \rangle - \min_{\mathbf{w} \in \mathcal{K}} \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w} \rangle \right) \\ &\stackrel{(3)}{=} \frac{1}{T} \left( \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w}_t \rangle - \min_{\mathbf{w} \in \mathcal{K}} \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w} \rangle \right). \end{aligned} \quad (5)$$

Here, Equality (1) follows from Equation (4), Inequality (2) holds because  $\mathbf{p}(\mathbf{x}_t, \mathbf{y}_t) \in H_{\mathbf{w}_t}$ , and Equality (3) holds from our definition of  $\mathbf{l}_t$ . ■

As a corollary, since for any  $\mathbf{x} \in \mathbb{R}^d$  and  $S \subseteq \mathbb{R}^d$ , the  $\ell_\infty$  distance is always less than equal to the  $\ell_2$  distance, i.e.,  $d_\infty(\mathbf{x}, S) \leq d_2(\mathbf{x}, S)$ , we obtain

$$d_\infty\left(\frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S\right) \leq d_2\left(\frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S\right) \leq \frac{1}{T} \left( \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w}_t \rangle - \min_{\mathbf{w} \in \mathcal{K}} \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w} \rangle \right).$$

*OLO regret upper-bound with Follow-the-Regularized-Leader algorithm.* To obtain the upper-bound on the regret of an OLO problems in terms of the  $\ell_\infty$  norm of its losses, we apply the Follow-the-Regularized-Leader algorithm with a  $\mu$ -strongly convex<sup>23</sup> regularizer with respect to the  $\ell_1$  norm. We use a regularizer with respect to the  $\ell_1$  norm, the dual of the  $\ell_\infty$  norm, because of the bound structure of the algorithm as stated in [Lemma 2](#). We elaborate further in the following lemmas.

<sup>22</sup>  $S^\circ$  is the polar cone of  $S$ , i.e.  $S^\circ := \{\mathbf{s} \in \mathbb{R}^d : \langle \mathbf{s}, \mathbf{x} \rangle \leq 0 \text{ for all } \mathbf{x} \in S\}$ , and  $B_2(1)$  is a Euclidian ball with radius 1, i.e.  $B_2(1) = \{\mathbf{w} \in \mathbb{R}^d : \|\mathbf{w}\|_2 \leq 1\}$ .

<sup>23</sup> A  $\mu$ -strongly convex function  $f$  with respect to the  $\ell_q$  norm is a differentiable function that satisfies  $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla(f(\mathbf{x})^T)(\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_q^2$  for some  $\mu > 0$ .

LEMMA 2. (*Shalev-Shwartz et al. (2012)*) Consider an OLO problem on a convex and compact decision space  $\mathcal{K} \subseteq \mathbb{R}^d$ . Applying Follow-the-Regularized-Leader algorithm with a regularizer  $R$ , where  $R: \mathbb{R}^d \rightarrow \mathbb{R}$  is a  $\mu$ -strongly convex function with respect to some norm  $\|\cdot\|$  for  $\mu > 0$ , implies

$$\frac{1}{T} \left( \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w}_t \rangle - \min_{\mathbf{w} \in \mathcal{K}} \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w} \rangle \right) \leq O(CB^{1/2}\mu^{-1/2}T^{-1/2}),$$

where  $B > 0$  upper bounds the function  $R$ ,  $C > 0$  upper bounds the dual norm of the loss vector  $\|\mathbf{l}\|_*$ <sup>24</sup>, and  $T$  is the number of rounds.

LEMMA 3. (*Shalev-Shwartz (2007)*) For  $q \in (1, 2)$ , the function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  defined as  $f(\mathbf{x}) = \frac{1}{2(q-1)} \|\mathbf{x}\|_q^2$  is strongly convex with respect to the  $\ell_q$  norm over  $\mathbb{R}^d$ . Recall that the  $\ell_q$  norm is defined as  $\|\mathbf{x}\|_q = (x_1^q + x_2^q + \dots + x_d^q)^{1/q}$  for  $\mathbf{x} \in \mathbb{R}^d$ .

To get a bound from Lemma 2 that depends on the upper-bound of the  $\ell_\infty$  norm of the loss vectors, we want a regularizer  $R$  that is  $\mu$ -strongly convex w.r.t the  $\ell_1$  norm for some  $\mu > 0$  (to be determined later). However, the function from Lemma 3 does not work for  $q = 1$ . To solve this, we set  $q$  to be greater than 1,  $q = \frac{\log(d)}{\log(d)-1}$  in particular, then bound the  $\ell_q$  norm from above using the  $\ell_1$  norm. Specifically, setting  $R(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_q^2$  and  $q = \frac{\log(d)}{\log(d)-1}$  we have

$$\begin{aligned} R(\mathbf{x}) &= (q-1) \cdot \frac{1}{2(q-1)} \|\mathbf{x}\|_q^2 \\ &\stackrel{(1)}{\geq} (q-1) \frac{1}{2(q-1)} \|\mathbf{x}\|_q^2 + (q-1) \nabla \left( \frac{1}{2(q-1)} \|\mathbf{x}\|_q^2 \right)^T (\mathbf{y} - \mathbf{x}) + (q-1) \cdot \frac{1}{2} \cdot \|\mathbf{y} - \mathbf{x}\|_q^2 \\ &\stackrel{(2)}{\geq} (q-1) \frac{1}{2(q-1)} \|\mathbf{x}\|_q^2 + (q-1) \nabla \left( \frac{1}{2(q-1)} \|\mathbf{x}\|_q^2 \right)^T (\mathbf{y} - \mathbf{x}) + (q-1) \cdot \frac{1}{2} \cdot \frac{\|\mathbf{y} - \mathbf{x}\|_1^2}{3} \\ &= R(\mathbf{x}) + R(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{(q-1)/3}{2} \|\mathbf{y} - \mathbf{x}\|_1^2, \end{aligned}$$

where  $\mu = \frac{q-1}{3} = \frac{\frac{\log(d)}{\log(d)-1}-1}{3} = \frac{1}{3\log(d)}$ . So, the function  $R$  is  $\frac{1}{3\log(d)}$ -strongly convex with respect to the  $\ell_1$  norm. Furthermore, Inequality (1) follows from Lemma 3 and Inequality (2) holds because  $\|\mathbf{w}\|_1/3 \leq \|\mathbf{w}\|_q$  for any  $\mathbf{w} \in \mathbb{R}^d$ .

Consequently, by constructing an OLO problem with  $\mathcal{K} = S^\circ \cap B_2(1)$  and  $\mathbf{l}_t = -p(\mathbf{x}_t, \mathbf{y}_t)$  on each round, applying the Follow-the-Regularized-Leader algorithm with regularizer  $R(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_q^2$  to the OLO problem, and converting  $\mathbf{w}_t$  to  $\mathbf{x}_t$  in each round, we obtain

$$d_\infty \left( \frac{1}{T} \sum_{t=1}^T p(\mathbf{x}_t, \mathbf{y}_t), S \right) \leq \frac{1}{T} \left( \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w}_t \rangle - \min_{\mathbf{w} \in \mathcal{K}} \sum_{t=1}^T \langle \mathbf{l}_t, \mathbf{w} \rangle \right) \leq O(D(\mathbf{p}) \log(d)^{1/2} T^{-1/2}).$$

The last inequality follows from applying Lemma 2 to the OLO problem corresponding to the Blackwell game with  $B = 1$  and  $C = D(\mathbf{p})$ . Notice that for any  $\mathbf{w} \in \mathcal{K}$ ,  $\|\mathbf{w}\|_2 \leq 1 = B$  because we set  $\mathcal{K} = S^\circ \cap B_2(1)$  in Lemma 1. Furthermore, since we set  $\mathbf{l}_t = p(\mathbf{x}_t, \mathbf{y}_t)$ , we have  $\|\mathbf{l}_t\|_\infty = \|p(\mathbf{x}_t, \mathbf{y}_t)\|_\infty \leq D(\mathbf{p})$  by definition. ■

<sup>24</sup> If  $\|\cdot\|$  is a norm in  $\mathbb{R}^d$ , the dual norm  $\|\cdot\|_*$  of  $\|\cdot\|$  is defined as  $\|\mathbf{w}\|_* = \sup\{\mathbf{w}^T \mathbf{x} \mid \|\mathbf{x}\| \leq 1\}$ .

## Appendix C: Proofs and Remarks of Section 3.1

EXAMPLE 2 (NON-ROBUST GREEDY ALGORITHM). In the shortest path tree problem, we are given an undirected graph  $G = (V, E)$  along with a root node  $u$  and edge weights  $\{w_{uv}\}_{(u,v) \in E}$ . We want to compute a spanning tree of  $G$  such that for all vertices  $v \in V$ , the distance to the root in the tree,  $\text{dist}_T(u, v)$ , equals the distance to the root in the original graph,  $\text{dist}_G(u, v)$ . This problem can be solved by a greedy algorithm which runs Dijkstra’s algorithm from  $u$  and then for each node  $v \neq u$  chooses a parent  $p \in \text{neighborhood}(v)$  with the smallest  $w_{vp} + \text{dist}_G(p, u)$ . Suppose that we want to solve the online problem where the  $G$  and  $u$  are fixed over all rounds but the edge weights are chosen by an adversary.

This can be translated into the language of our meta-algorithm as follows. The feasible region is to choose a parent for every non-root vertex ( $\mathcal{C} = \prod_{v \in V \setminus \{u\}} \text{neighborhood}(v)$ ). The adversary’s function space is to choose (bounded) weights ( $\mathcal{F} \cong (0, 1]^E$ ), and the cost of a chosen set of edges that we aim to minimize is the average distance from a random vertex to  $u$ . For each of our  $|V|$  subroutines, the parameter space is to choose a distribution for the parent vertex ( $\Theta = \Delta(\text{neighborhood}(v))$ )<sup>25</sup>. The (one-dimensional) payoff vector is the shortest path from  $v$  to  $u$  through that parent  $p$  ( $\text{PAYOFF}(\boldsymbol{\theta}, \mathbf{z}, \{w_{uv}\}) = \mathbb{E}_{p \sim \boldsymbol{\theta}} [w_{vp} + \text{dist}_G(p, u)]$ ), where  $\boldsymbol{\theta}$  is the probability of choosing a vertex among  $v$ ’s neighbors as parent.

Managing to perfectly minimize the one-dimensional payoff vector at each iteration results in a shortest path tree and therefore the best possible objective value. However, if the local choices deviate from their optimal values, then we can create cycles which result in infinite objective value.

For example, consider the clique on  $V = \{1, 2, 3\}$ , where we want a shortest path tree to the root node  $u = 1$ . When the weights are  $w_{12} = 0.25, w_{13} = 1.0, w_{23} = 0.5$ , then it would be best for node three to first take edge  $(2, 3)$ . If we swap the role of nodes two and three,  $w_{12} = 1.0, w_{13} = 0.25, w_{23}$ , then it would be best for node two to first take edge  $(2, 3)$ . When our subroutines for nodes two and three make simultaneous decisions without actually seeing the input, they could easily both choose edge  $(2, 3)$ , yielding an invalid shortest path tree and making it impossible to get from either node to the root. This global issue can’t be expressed as local utilities, so the algorithm is not robust in the sense that is needed to apply our framework.

## Appendix D: Proofs and Remarks of Section 5

### D.1. Proof of Theorem 3

In this section, we complete the proof of Theorem 3, which is restated below for convenience.

THEOREM 3. A closed convex set  $S$  is  $O\left(D(\mathbf{p})^{1/3} D(\hat{\mathbf{p}})^{2/3} (\log d)^{1/3} T^{-1/3}\right)$ -bandit-approachable in the bandit Blackwell sequential game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$  if and only if  $S$  is response-satisfiable in the Blackwell game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$ . In particular, when  $S$  is response satisfiable, the online algorithm AlgBB (Algorithm 3) achieves this approachability bound in polynomial time, given access to a separation oracle for  $S$ .

*Proof.* The only if direction is proved in the sketch. To prove the if direction and second part of Theorem 3, we propose an algorithm AlgBB that is parameterized by an exploration probability  $q \in (0, 1]$  (Algorithm 3). We later choose  $q$  to be  $D(\mathbf{p})^{-2/3} D(\hat{\mathbf{p}})^{2/3} (\log d)^{1/3} T^{-1/3}$  to balance terms in our regret upper-bound. In each round  $t$ , this algorithm outputs a move  $\mathbf{x}_t \in \mathcal{Y}$  as well as whether to explore  $\pi_t \in \{\text{YES}, \text{NO}\}$ ,

<sup>25</sup> Our framework is amenable to the parameter space depending on the iteration  $i$ .

and then receives an unbiased estimate of the resulting payoff  $\hat{\mathbf{p}}(\mathbf{x}_t, \mathbf{y}_t)$  based on both players' actions if it picks to explore. It also maintains a (full-information) Blackwell algorithm AlgB. In each round  $t = 1, 2, \dots, T$ , our algorithm follows the last suggested action by AlgB to generate a move  $\mathbf{x}_t$ . Note that this move will be exactly the same as the previous round, if the algorithm chose to not explore in the previous round. Our algorithm then decides to either explore with probability  $q$  or not explore with probability  $1 - q$ . If it explores, then it receives an unbiased estimator,  $\hat{\mathbf{p}}$  of the current vector payoff  $\mathbf{p}(\mathbf{x}_t, \mathbf{y}_t)$ , and passes a scaled version  $\hat{\mathbf{p}}/q$  on to AlgB. If it does not explore, then it rewinds the state of AlgB to the beginning of the current round. Our goal here is to show that under algorithm AlgBB,  $d_\infty\left(\frac{1}{T}\sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S\right)$  plus exploring penalty term  $\mathbb{E}\left[\frac{1}{T}D(\mathbf{p}) \cdot (\# \text{ explore})\right]$  is  $O\left(D(\mathbf{p})^{1/3}D(\hat{\mathbf{p}})^{2/3}(\log d)^{1/3}T^{-1/3}\right)$ .

We start with bounding the first term  $d_\infty\left(\frac{1}{T}\sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S\right)$  as a function of  $\hat{\Pi}_T$ , which is the time-averaged rescaled estimated payoffs from the rounds that we explore in  $1, 2, \dots, T$ :

$$\hat{\Pi}_T \triangleq \frac{1}{T} \sum_{t=1}^T \frac{1}{q} \hat{\mathbf{p}}(\mathbf{x}_t, \mathbf{y}_t) \mathbb{1}[\text{explore in round } t].$$

Specifically, we have

$$\begin{aligned} d_\infty\left(\frac{1}{T}\sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S\right) &= d_\infty\left(\frac{1}{T}\sum_{t=1}^T \mathbb{E}\left[\frac{1}{q}\hat{\mathbf{p}}(\mathbf{x}_t, \mathbf{y}_t) \mathbb{1}[\text{explore in round } t]\right], S\right) \\ &\leq \mathbb{E}\left[d_\infty\left(\hat{\Pi}_T, S\right)\right], \end{aligned}$$

where the equality follows because  $\mathbb{E}\left[\frac{1}{q}\hat{\mathbf{p}}(\mathbf{x}_t, \mathbf{y}_t) \mathbb{1}[\text{explore in round } t]\right] = \frac{q}{q}\mathbb{E}[\hat{\mathbf{p}}(\mathbf{x}_t, \mathbf{y}_t)] = \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t)$  as  $\hat{\mathbf{p}}$  is an unbiased estimator for  $\mathbf{p}$ , and the inequality is obtained by applying Jensen's inequality to the convex  $\ell_\infty$  distance function. We next show that if we explore with probability  $q$ ,  $\mathbb{E}\left[d_\infty\left(\hat{\Pi}_T, S\right)\right] \leq O(D(\hat{\mathbf{p}})\log(d)^{1/2}(qT)^{-1/2})$ . Then, observe that the exploring penalty term  $\mathbb{E}\left[\frac{1}{T}D(\mathbf{p}) \cdot (\# \text{ explore})\right]$  equals  $D(\mathbf{p})q$ . Our choice of exploring probability  $q = D(\mathbf{p})^{-2/3}D(\hat{\mathbf{p}})^{2/3}\log(d)^{1/3}T^{-1/3}$  makes the two terms equal to  $O(D(\mathbf{p})^{1/3}D(\hat{\mathbf{p}})^{2/3}(\log d)^{1/3}T^{-1/3})$ , and gives us the desired bound.

To see why  $\mathbb{E}\left[d_\infty\left(\hat{\Pi}_T, S\right)\right] \leq O(D(\hat{\mathbf{p}})\log(d)^{1/2}(qT)^{-1/2})$  when we explore with probability  $q$ , let  $M$  be a random variable equal to the number of rounds we explore and  $(\tau_1, \tau_2, \dots, \tau_M)$  be the rounds that we explore. Note that  $M \sim \text{Binomial}(T, q)$ . By applying the law of total expectation, we have:

$$\mathbb{E}\left[d_\infty\left(\hat{\Pi}_T, S\right)\right] = \sum_{m=0}^T \mathbb{E}\left[d_\infty\left(\hat{\Pi}_T, S\right) \mid M = m\right] \Pr[M = m].$$

We provide an upper bound on each term in the above summation separately. First, we handle the  $M = m = 0$  case by noting  $\hat{\Pi}_T = \mathbf{0}$ , hence the distance from  $S$  is bounded by  $D(\hat{\mathbf{p}})$  in this case. Moreover, this event occurs with probability  $(1 - q)^T$ . See that

$$(1 - q)^T = (1 - q)^{q^{-1} \cdot q \cdot T} \leq (1/e)^{qT} \leq O((qT)^{-1/2}),$$

and therefore  $\mathbb{E}\left[d_\infty\left(\hat{\Pi}_T, S\right) \mid M = 0\right] \Pr[M = 0] \leq O\left(D(\hat{\mathbf{p}})\log(d)^{1/2}(qT)^{-1/2}\right)$ .

Now fix some  $M = m \neq 0$ . Assuming that  $S$  is a cone (we can always lift the convex set  $S$  to a cone in one dimension higher as shown in [Abernethy et al. \(2011\)](#)), our full-information Blackwell algorithm AlgB, who receives "fake payoffs"  $\{\hat{\mathbf{p}}(\mathbf{x}_{\tau_i}, \mathbf{y}_{\tau_i})\}_{i=1}^m$  with a diameter of  $\frac{1}{q}D(\hat{\mathbf{p}})$ , guarantees that:

$$\mathbb{E}\left[d_\infty\left(\hat{\Pi}_T, S\right) \mid M = m\right] = \mathbb{E}\left[\frac{M}{T} \cdot d_\infty\left(\frac{T}{M}\hat{\Pi}_T, S\right) \mid M = m\right]$$

$$\begin{aligned}
&= \frac{m}{T} \mathbb{E} \left[ d_\infty \left( \frac{1}{M} \sum_{i=1}^M \frac{1}{q} \hat{\mathbf{p}}(\mathbf{x}_{\tau_i}, \mathbf{y}_{\tau_i}), S \right) \middle| M = m \right] \\
&\leq \frac{m}{T} \cdot O \left( \frac{1}{q} D(\hat{\mathbf{p}}) \log(d)^{1/2} m^{-1/2} \right) \\
&= O \left( \frac{1}{qT} D(\hat{\mathbf{p}}) \log(d)^{1/2} m^{1/2} \right), \tag{6}
\end{aligned}$$

where the expectation is taken w.r.t. the randomness in  $\hat{\mathbf{p}}$  and the inequality holds because  $S$  is response-satisfiable in the Blackwell game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  and  $\hat{\mathbf{p}}$  is an unbiased estimator of  $\mathbf{p}$ . To be more clear why the above inequality holds, note that set  $S$  is response-satisfiable in the Blackwell game  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$ , and is not necessarily response-satisfiable if we replace  $\mathbf{p}$  with  $\hat{\mathbf{p}}$ . However, by (i) following exactly the same steps as in proof of Theorem 1 (Section B.2 in the appendix) to reduce Blackwell approachability to online linear optimization for rounds  $\{\tau_i\}_{i=1}^M$ , (ii) plugging  $\hat{\mathbf{p}}$  as the vector payoff of each round and using  $\mathbf{l}_i = -\hat{\mathbf{p}}(\mathbf{x}_{\tau_i}, \mathbf{y}_{\tau_i})$  as the loss function in the online linear optimization, and then (iii) using the fact that  $\hat{\mathbf{p}}$  is an unbiased estimator for  $\mathbf{p}$  and  $S$  is response-satisfiable w.r.t. payoffs  $\mathbf{p}$ , we can obtain exactly the same approachability bound in expectation as if  $S$  was response-satisfiable w.r.t. payoffs  $\hat{\mathbf{p}}$ . To see this, consider the chain of inequalities (5) in the proof of Theorem 1 in Section B.2, tailored to our problem, and take an expectation w.r.t. the randomness in  $\hat{\mathbf{p}}$ . We have:

$$\begin{aligned}
\mathbb{E} \left[ d_\infty \left( \frac{1}{M} \sum_{i=1}^M \hat{\mathbf{p}}(\mathbf{x}_{\tau_i}, \mathbf{y}_{\tau_i}), S \right) \middle| \{\tau_i\}_{i=1}^M \right] &\leq \mathbb{E} \left[ \max_{\mathbf{w} \in \mathcal{K}} \left\langle \frac{1}{M} \sum_{i=1}^M \hat{\mathbf{p}}(\mathbf{x}_{\tau_i}, \mathbf{y}_{\tau_i}), \mathbf{w} \right\rangle \middle| \{\tau_i\}_{i=1}^M \right] \\
&= \mathbb{E} \left[ \frac{1}{M} \max_{\mathbf{w} \in \mathcal{K}} \left( - \sum_{i=1}^M \langle \mathbf{l}_i, \mathbf{w} \rangle \right) \middle| \{\tau_i\}_{i=1}^M \right] \\
&\stackrel{(2)}{\leq} \frac{1}{M} \left( \sum_{i=1}^M \langle -\mathbf{p}(\mathbf{x}_{\tau_i}, \mathbf{y}_{\tau_i}), \mathbf{w}_i \rangle - \mathbb{E} \left[ \min_{\mathbf{w} \in \mathcal{K}} \sum_{i=1}^M \langle \mathbf{l}_i, \mathbf{w} \rangle \middle| \{\tau_i\}_{i=1}^M \right] \right) \\
&\stackrel{(3)}{=} \frac{1}{M} \mathbb{E} \left[ \sum_{i=1}^M \langle \mathbf{l}_i, \mathbf{w}_i \rangle - \min_{\mathbf{w} \in \mathcal{K}} \sum_{i=1}^M \langle \mathbf{l}_i, \mathbf{w} \rangle \middle| \{\tau_i\}_{i=1}^M \right].
\end{aligned}$$

This time, Inequality (2) holds as before, because  $S$  is response-satisfiable w.r.t. payoffs  $\mathbf{p}$  (and hence half-space satisfiable when using  $\mathbf{w}_t$  as the normal of the half-space), but Equality (3) holds because:

$$-\langle \mathbf{p}(\mathbf{x}_{\tau_i}, \mathbf{y}_{\tau_i}), \mathbf{w}_i \rangle = -\mathbb{E} [\langle \hat{\mathbf{p}}(\mathbf{x}_{\tau_i}, \mathbf{y}_{\tau_i}), \mathbf{w}_i \rangle | \{\tau_i\}_{i=1}^M] = -\mathbb{E} [\langle \mathbf{l}_i, \mathbf{w}_i \rangle | \{\tau_i\}_{i=1}^M]$$

Note that expectation is conditioned on  $\{\tau_i\}_{i=1}^M$ , but only we use a universal upper-bound on the last term (regret of online linear optimization) that is a function of  $M$ , so we can change the conditioning on only  $M$ .

We now use the bound in (6) for  $q = D(\mathbf{p})^{-2/3} D(\hat{\mathbf{p}})^{2/3} \log(d)^{1/3} T^{-1/3}$ , and then use Jensen's inequality applied to the (concave) square-root function.

$$\begin{aligned}
\mathbb{E} \left[ d_\infty \left( \hat{\Pi}_T, S \right) + \frac{1}{T} D(\mathbf{p}) \cdot (\# \text{ explore}) \right] &= \mathbb{E}_{m \sim \text{Binomial}(T, q)} \left[ \mathbb{E} \left[ d_\infty \left( \hat{\Pi}_T, S \right) \middle| M = m \right] \right] + O(D(\mathbf{p})q) \\
&\leq \mathbb{E}_{m \sim \text{Binomial}(T, q)} \left[ O \left( \frac{1}{qT} D(\hat{\mathbf{p}}) \log(d)^{1/2} m^{1/2} \right) \right] + O(D(\mathbf{p})q) \\
&\leq O \left( \frac{1}{qT} D(\hat{\mathbf{p}}) \log(d)^{1/2} (Tq)^{1/2} \right) + O(D(\mathbf{p})q) \\
&= O(D(\hat{\mathbf{p}}) \log(d)^{1/2} (qT)^{-1/2}) + O(D(\mathbf{p})q) \\
&= O \left( D(\mathbf{p})^{1/3} D(\hat{\mathbf{p}})^{2/3} (\log d)^{1/3} T^{-1/3} \right).
\end{aligned}$$

The last inequality is the desired result. ■

## D.2. Proof of Theorem 4

*Proof.* The function  $\hat{\mathbf{p}}$  is an unbiased estimator for  $\mathbf{p}$  (due to the bandit Blackwell reducibility), so  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$  is a valid instance of the bandit Blackwell sequential game. Moreover, our target set  $S$  is the  $d_{\text{payoff}}$ -dimensional positive orthant. Therefore, there exists a polynomial-time separation oracle for set  $S$ . Set  $S$  is also response-satisfiable (due to bandit Blackwell reducibility). Thus, there exists a polynomial time online algorithm AlgBB that guarantees the bandit approachability upper-bound  $O\left(D(\mathbf{p})^{1/3}D(\hat{\mathbf{p}})^{2/3}(\log(d_{\text{payoff}}))^{1/3}T^{2/3}\right)$ , established in Theorem 3, for each of the bandit Blackwell instances corresponding to the  $N$  different subproblems.

Consider a subproblem  $i \in [N]$ . Note that AlgBB<sup>(i)</sup> is not invoked in all rounds  $[T]$ , but rather a subset  $\mathcal{T}_i \subseteq [T]$  depending on when its fellow Blackwell bandit algorithms, i.e., AlgBB<sup>(i')</sup>,  $i' \in [i-1]$ , decide to explore. Note that  $\mathcal{T}_i$  is a random set, and only depends on realizations of binary signals  $\{\pi_t^{(i')}\}_{i' \in [i-1], t \in [T]}$ . Fix a particular realization of set  $\mathcal{T}_i$ . By using the upper-bound in Theorem 3 for each of the terms in the LHS of the bound (i.e., the distance of the average payoff vector from set  $S$  and expected number of explorations) separately, we have

$$d_\infty \left( \frac{1}{|\mathcal{T}_i|} \sum_{t \in \mathcal{T}_i} \mathbf{p} \left( \boldsymbol{\theta}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t) \right), S \right) \leq O \left( D(\mathbf{p})^{1/3} D(\hat{\mathbf{p}})^{2/3} (\log(d_{\text{payoff}}))^{1/3} |\mathcal{T}_i|^{-1/3} \right).$$

Moreover, let  $\mathcal{T}_i^+$ ,  $\mathcal{T}_i^-$ ,  $M_i$  be the rounds where AlgBB<sup>(i)</sup> explores, the rounds where AlgBB<sup>(i)</sup> exploits, and the number of rounds that AlgBB<sup>(i)</sup> explores respectively, i.e.  $M_i = |\mathcal{T}_i^+|$ . Then, by our choice of  $q = D(\mathbf{p})^{-2/3}D(\hat{\mathbf{p}})^{2/3}(\log d)^{1/3}T^{-1/3}$ , we have

$$\begin{aligned} \mathbb{E} \left[ \frac{1}{|\mathcal{T}_i|} M_i \middle| \mathcal{T}_i \right] &= D(\mathbf{p})^{-2/3} D(\hat{\mathbf{p}})^{2/3} (\log(d_{\text{payoff}}))^{1/3} |T|^{-1/3} \\ &\leq D(\mathbf{p})^{-2/3} D(\hat{\mathbf{p}})^{2/3} (\log(d_{\text{payoff}}))^{1/3} |\mathcal{T}_i|^{-1/3}, \end{aligned}$$

Let  $\mathcal{T}^-$  be the set of rounds where no AlgBB<sup>(i)</sup> explored and  $\mathcal{T}^+$  be the set of rounds where some AlgBB<sup>(i)</sup> explored. Note also that  $\mathcal{T}^- \subseteq \mathcal{T}_i$ , simply because if no algorithm explores, AlgBB<sup>(i)</sup> will be invoked. Notice that for the rounds where AlgBB<sup>(i)</sup> is invoked and not exploring, but there exists some subroutine  $j > i$  that explores, no feedback was given to AlgBB<sup>(i)</sup>, so we can think of it as if AlgBB<sup>(i)</sup> is not being invoked. Hence, combining this with the fact that the set  $S$  is the positive orthant, we have:

$$\forall j \in [n]: \left[ \sum_{t \in (\mathcal{T}_i \setminus \mathcal{T}^+) \cup \mathcal{T}_i^+} \mathbf{p} \left( \boldsymbol{\theta}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t) \right) \right]_j \geq -O \left( D(\mathbf{p})^{1/3} D(\hat{\mathbf{p}})^{2/3} (\log(d_{\text{payoff}}))^{1/3} |(\mathcal{T}_i \setminus \mathcal{T}^+) \cup \mathcal{T}_i^+|^{2/3} \right).$$

Furthermore, due to Blackwell reducibility, which is a precondition for bandit Blackwell reducibility (see Definitions 7 and 9),  $\text{PAYOFF} \left( \boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t \right) = \mathbf{p} \left( \boldsymbol{\theta}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t) \right)$ . Thus,

$$\forall j \in [n]: \left[ \sum_{t \in (\mathcal{T}_i \setminus \mathcal{T}^+) \cup \mathcal{T}_i^+} \text{PAYOFF}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) \right]_j \geq -O \left( D(\mathbf{p})^{1/3} D(\hat{\mathbf{p}})^{2/3} (\log(d_{\text{payoff}}))^{1/3} |(\mathcal{T}_i \setminus \mathcal{T}^+) \cup \mathcal{T}_i^+|^{2/3} \right). \quad (7)$$

For the rounds where no AlgBB<sup>(i)</sup> explore, for any  $j \in [n]$ , we have

$$\begin{aligned} \mathbb{E} \left[ \sum_{t \in \mathcal{T}^-} \text{PAYOFF}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) \middle| \mathcal{T}_i \right]_j &= \mathbb{E} \left[ \sum_{t \in (\mathcal{T}_i \setminus \mathcal{T}^+) \cup \mathcal{T}_i^+} \text{PAYOFF}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) \middle| \mathcal{T}_i \right]_j \\ &\quad - \mathbb{E} \left[ \sum_{t \in \mathcal{T}_i^+} \text{PAYOFF}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) \middle| \mathcal{T}_i \right]_j \\ &\geq -O \left( D(\mathbf{p})^{1/3} D(\hat{\mathbf{p}})^{2/3} (\log(d_{\text{payoff}}))^{1/3} |(\mathcal{T}_i \setminus \mathcal{T}^+) \cup \mathcal{T}_i^+|^{2/3} \right) - D(\mathbf{p}) \mathbb{E}[M_i | \mathcal{T}_i], \end{aligned}$$

where the expectation is with respect to  $\mathbf{z}_t^{(i-1)}$ ,  $t \in \mathcal{T}_i$ . The equality holds because  $\mathcal{T}^- \cup \mathcal{T}_i^+ = (\mathcal{T}_i \setminus \mathcal{T}^+) \cup \mathcal{T}_i^+$  and  $\mathcal{T}^- \cap \mathcal{T}_i^+ = \emptyset$ . The inequality follows from Equation (7) and the fact that  $M_i = |\mathcal{T}_i^+|$  and for any  $i \in [N]$  and  $t \in [T]$ ,  $\text{PAYOFF}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) \leq D(\mathbf{p})$ . By considering the fact that

$$\mathbb{E}[M_i | \mathcal{T}_i] = D(\mathbf{p})^{-2/3} D(\hat{\mathbf{p}})^{2/3} (\log d)^{1/3} T^{-1/3} |\mathcal{T}_i| \leq D(\mathbf{p})^{-2/3} D(\hat{\mathbf{p}})^{2/3} (\log(d_{\text{payoff}}))^{1/3} |\mathcal{T}_i|^{2/3}$$

from our choice of the probability of exploring,  $q$ , in Theorem 3,  $|\mathcal{T}_i| \leq T$  and  $|(\mathcal{T}_i \setminus \mathcal{T}^+) \cup \mathcal{T}_i^+| \leq T$ , we have:

$$\forall j \in [n]: \mathbb{E} \left[ \sum_{t \in \mathcal{T}^-} \text{PAYOFF}(\boldsymbol{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) \right]_j \geq -O \left( D(\mathbf{p})^{1/3} D(\hat{\mathbf{p}})^{2/3} (\log(d_{\text{payoff}}))^{1/3} T^{2/3} \right). \quad (8)$$

Because the offline algorithm OFFLINE-IG( $\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta$ ) (Algorithm 1) is an extended  $(\gamma, \delta)$ -robust approximation, by focusing on rounds in  $\mathcal{T}^-$  and applying Inequality (8), together with linearity of expectation, we have:

$$\mathbb{E} \left[ \sum_{t \in \mathcal{T}^-} f_t(\mathbf{z}_t) \right] \geq \gamma \cdot \mathbb{E} \left[ \sum_{t \in \mathcal{T}^-} f_t(\mathbf{z}^*) \right] - O \left( D(\mathbf{p})^{1/3} D(\hat{\mathbf{p}})^{2/3} N \delta (\log(d_{\text{payoff}}))^{1/3} T^{2/3} \right),$$

where  $\mathbf{z}^* = \arg \max_{\mathbf{z} \in \mathcal{C}} \sum_{t=1}^T f_t(\mathbf{z})$  is the optimal in-hindsight solution.

Finally, note that BANDIT-IG( $\mathcal{C}, \mathcal{F}, \mathcal{D}, \Theta, \text{AlgBB}$ ) does not explore too often in total among its subproblems. More precisely,

$$\mathbb{E}[|\mathcal{T}^+|] = \mathbb{E} \left[ \sum_{i=1}^N M_i \right] \leq \sum_{i=1}^N O \left( (\log(d_{\text{payoff}}))^{1/3} \mathbb{E}[\mathcal{T}_i^{2/3}] \right) \leq O \left( N (\log(d_{\text{payoff}}))^{1/3} T^{2/3} \right).$$

Noting the fact that the functions  $f_t$  have output value at most 1, for the remaining rounds  $\mathcal{T}^+$  we have:

$$\mathbb{E} \left[ \sum_{t \in \mathcal{T}^+} f_t(\mathbf{z}_t) \right] \geq \gamma \cdot \mathbb{E} \left[ \sum_{t \in \mathcal{T}^+} f_t(\mathbf{z}^*) \right] - O \left( N (\log(d_{\text{payoff}}))^{1/3} T^{2/3} \right). \quad (9)$$

Combining the two types of bounds in rounds  $\mathcal{T}^-$  and  $\mathcal{T}^+$  yields the desired claim. ■

### D.3. Bandit Blackwell Regret Lowerbound

In this section, we show that in a Bandit Blackwell sequential game,  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$ , the distance from the time-averaged payoff to the target set  $S$  plus the time-averaged exploring penalty of any prediction strategy must be at least  $\Omega(\underline{D}T^{-1/3})$ , where  $\underline{D} = \min\{D(\mathbf{p}), D(\hat{\mathbf{p}})\}$ . Put differently, we show that the performance bound proved in Theorem 3 is unimprovable with respect to  $T$  (the number of rounds), i.e., no other strategies can have a better performance for all problems.

**THEOREM 7.** *In a bandit Blackwell sequential game,  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}, \hat{\mathbf{p}})$ , there exists an adversary's strategy such that for every player 1's strategy, the resulting sequence of actions satisfy:*

$$d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) + \mathbb{E} \left[ \frac{1}{T} D(\mathbf{p}) \cdot (\# \text{ explore}) \right] \geq \Omega(\underline{D} T^{-1/3}).$$

where  $(\# \text{ explore})$  is the number of exploration rounds and  $\underline{D} = \min\{D(\mathbf{p}), D(\hat{\mathbf{p}})\}$ .

*Proof of Theorem 7.* Let  $M$  be a random variable equal to the number of rounds the player explores. We first show in that if the number of rounds that the player explores at is at most  $M$ , then there exists a Bandit Blackwell instance: an adversary's action  $(\mathbf{y}_1, \dots, \mathbf{y}_T)$ , a convex closed set  $S$ , and an affine payoff  $\mathbf{p}$  together with an unbiased estimator function  $\hat{\mathbf{p}}$  such that:

$$\mathbb{E} \left[ d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) \middle| M \right] \geq \Omega \left( \frac{D}{\sqrt{M}} \right)$$

for any player's strategies  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$ , where the expectation is taken over the randomness in the adversary's strategy. We later show this statement in Lemma 4. For now, we assume that the statement is true. Since the Bandit Blackwell total regret defined in Definition 8 includes another term for the cost of exploring, the total regret conditioned on  $M$  is

$$\begin{aligned} \mathbb{E} \left[ d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) \middle| M \right] + \mathbb{E} \left[ \frac{1}{T} D(\mathbf{p}) \cdot (\# \text{ explore}) \middle| M \right] &\geq \Omega \left( \frac{D}{\sqrt{M}} \right) + \Omega \left( \frac{DM}{T} \right) \\ &\stackrel{(1)}{\geq} \Omega(\underline{D} T^{-1/3}), \end{aligned}$$

where Inequality (2) follows from setting  $M = T^{2/3}$ ; notice that at  $M = T^{2/3}$ ,  $\frac{D}{\sqrt{M}} = \frac{DM}{T}$  and  $\Omega \left( \frac{D}{\sqrt{M}} \right) + \Omega \left( \frac{DM}{T} \right)$  is minimized. Again, the expectation here is with respect to the adversary's strategy. Now, taking another expectation with respect to  $M$ , we have

$$\begin{aligned} &\mathbb{E} \left[ d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) \right] + \mathbb{E} \left[ \frac{1}{T} D(\mathbf{p}) \cdot (\# \text{ explore}) \right] \\ &= \mathbb{E} \left[ \mathbb{E} \left[ d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) \middle| M \right] \right] + \mathbb{E} \left[ \mathbb{E} \left[ \frac{1}{T} D(\mathbf{p}) \cdot (\# \text{ explore}) \middle| M \right] \right] \geq \Omega(\underline{D} T^{-1/3}). \end{aligned}$$

This completes the proof. ■

We now prove the lower bound on the distance from the average payoff to  $S$  when the number of exploration rounds is  $M$ . As is common in proofs of lower bounds, we construct a random sequence of similar adversaries and show that with  $M$  rounds of explorations, it is impossible to distinguish the different types of adversaries without suffering a regret less than  $\frac{D}{\sqrt{M}}$ .

**LEMMA 4.** *In a Bandit Blackwell problem, if the number of exploration rounds is at most  $M$ , there exists an adversary's strategy  $(\mathbf{y}_1, \dots, \mathbf{y}_T)$ , a convex closed set  $S$ , and an affine payoff  $\mathbf{p}$  together with an unbiased estimator function  $\hat{\mathbf{p}}$  such that for any strategies  $(\mathbf{x}_1, \dots, \mathbf{x}_T)$ ,*

$$\mathbb{E} \left[ d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) \middle| M \right] \geq \Omega \left( \frac{D}{\sqrt{M}} \right),$$

where the expectation is taken with respect to the adversary's strategies.

*Proof of Lemma 4* We only prove our lower bound for a deterministic player. Note that any randomized strategy can be expressed as a randomization of deterministic strategies, and based on Yao's minimax principle (Yao (1977)), our lower bound still holds when we average them over several deterministic strategies according to some randomization. We refer to Cesa-Bianchi and Lugosi (2006) for details on deriving an identical lower bound for a randomized adversary from a deterministic adversary.

Recall that  $M$  is the random variable of the number of rounds that the player explores. Consider a fixed  $M$ . From this point onward (until the end of the proof), every probability and expectation are conditioned on  $M$ , we remove the dependence on  $M$  on the notations for simplicity. Let  $\mathcal{X} = \Delta([n])$ ,  $\mathcal{Y} = \{0, 1\}^n$ , the payoff function  $\mathbf{p}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y} \mathbf{1}_n - \mathbf{y}$ , and  $S$  be the non-positive orthant, i.e.  $S = \{\mathbf{s} \in \mathbb{R}^n | \mathbf{s}_j \leq 0 \ \forall j \in [n]\}$ . For deterministic strategies,  $\mathbf{x}$  must be equal to  $\mathbf{e}_z$  for some coordinate  $z$ , which happens when player 1 plays action  $z \in [n]$ . In that case,  $[\mathbf{p}(\mathbf{x}, \mathbf{y})]_j = [\mathbf{y}]_z - [\mathbf{y}]_j$  for all  $j \in [n]$ .

We now define the adversary's strategy. For each round  $t \in [T]$  and coordinate  $j \in [n]$ , let  $[\mathbf{y}_t]_j$  be Bernoulli random variables whose joint distribution are defined as follows. We first pick a random variable  $\zeta \sim \text{Uniform}\{1, 2, \dots, n\}$ . Then, given that  $\zeta = i$ ,  $[\mathbf{y}_1]_j, [\mathbf{y}_2]_j, \dots, [\mathbf{y}_T]_j$  are conditionally independent Bernoulli random variables with parameter  $(1 - \mu)/2$  if  $j \neq i$ , and  $(1 + \mu)/2$  if  $j = i$ , where  $\mu < 1/4$  (will be specified later). For analysis purposes, we define another move for the adversary, which we call the base move: all  $[\mathbf{y}_1]_j, [\mathbf{y}_2]_j, \dots, [\mathbf{y}_T]_j$  are conditionally independent Bernoulli variables with parameter  $(1 - \mu)/2$ . Suppose that this happens when  $\zeta = 0$  (just for ease of notations).

Let  $I_t$  be the player's action (in  $\{1, 2, \dots, n\}$ ) on round  $t$ , and  $\pi_t$  be the exploring indicator in round  $t$ :  $\pi_t = 1$  if we explore in round  $t$ , and 0 otherwise. Let  $\boldsymbol{\eta}_t = (\pi_1, \dots, \pi_t)$  be the history of exploration decisions up to round  $t$ . Since the player is deterministic,  $I_t$  is determined by  $(\mathbf{p}(\mathbf{x}_1, \mathbf{y}_1), \mathbf{p}(\mathbf{x}_2, \mathbf{y}_2), \dots, \mathbf{p}(\mathbf{x}_{t-1}, \mathbf{y}_{t-1}), \boldsymbol{\eta}_{t-1})$ . Also, let  $T_j = \sum_{t=1}^T \mathbb{1}[I_t = j]$  be the number of times action  $j$  is played in the first  $T$  rounds. We further define  $\mathbb{P}_j$  and  $\mathbb{E}_j$  as  $\mathbb{P}(\cdot | \zeta = j)$   $\mathbb{E}(\cdot | \zeta = j)$ , respectively. More rigorously, if  $\mathcal{A}$  is a  $\sigma$ -algebra generated by all possible outcomes of the game,  $\mathbb{P}_j$  is a measure on the  $\sigma$ -algebra  $\mathcal{A}$  and  $\mathbb{E}_j$  is an expectation taken with respect to the conditional probability  $\mathbb{P}_j$ , which solely depends on the adversary's move since we assume that player 1's strategy is deterministic.

Recall that for any  $j \in [n]$ , when  $\zeta = j$ , playing action  $j$  has the highest average reward than any other actions. Then, we have

$$\begin{aligned}
\mathbb{E}_j \left[ d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) \right] &\stackrel{(1)}{\geq} \max_{z \in [n]} \mathbb{E}_j \left[ \left[ \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t) \right]_z \right] \\
&= \frac{1}{T} \max_{z \in [n]} \mathbb{E}_j \left[ \sum_{t=1}^T ([\mathbf{y}_t]_z - [\mathbf{y}_t]_{I_t}) \right] \geq \frac{1}{T} \mathbb{E}_j \left[ \sum_{t=1}^T ([\mathbf{y}_t]_j - [\mathbf{y}_t]_{I_t}) \right] \\
&= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_j [[\mathbf{y}_t]_j - [\mathbf{y}_t]_{I_t}] \\
&\stackrel{(2)}{=} \frac{1}{T} \sum_{t=1}^T \mu \mathbb{E}_j [\mathbb{1}(I_t \neq j)] \\
&= \frac{1}{T} \mu \sum_{j' \neq j} \mathbb{E}_j [T_{j'}] = \frac{1}{T} \mu (T - \mathbb{E}_j [T_j]) \\
&= \left( \mu - \frac{\mu}{T} \mathbb{E}_j [T_j] \right).
\end{aligned}$$

Inequality (1) follows because  $S$  is the non-positive orthant. Equality (2) follows because  $\mathbb{E}_j[[\mathbf{y}_t]_{j'}]$  is  $(1+\mu)/2$  when  $j'=j$  and  $(1-\mu)/2$  otherwise, so the difference between  $\mathbb{E}_j[[\mathbf{y}_t]_j]$  and  $\mathbb{E}_j[[\mathbf{y}_t]_{I_t}]$  is  $\mu$  when  $I_t \neq j$  and 0 otherwise.

As for each  $j \in [n]$ , since the event  $\{\zeta = j\}$  happens with probability  $\frac{1}{n}$ , we have

$$\sup \mathbb{E} \left[ d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) \right] \geq \mu \left( 1 - \frac{1}{nT} \sum_j \mathbb{E}_j[T_j] \right), \quad (10)$$

where the expectation is taken with respect to the adversary's move and the supremum is taken over  $\zeta \in \{1, 2, \dots, n\}$  since the  $\zeta$  picked by the adversary in the beginning of the game determines his whole strategy. The proof now reduces to bounding  $\mathbb{E}_j[T_j]$  from above. We do this by comparing  $\mathbb{E}_j[T_j]$  with  $\mathbb{E}_0[T_j]$ . If player 1 chooses action  $i$  at round  $t$  and decides to explore, i.e.,  $I_t = i$  and  $\pi_t = 1$ , he then observes the payoff  $[\mathbf{y}_t]_i$ . Recall that  $\mathbf{y}_t$  is the random variable that represents adversary's move at round  $t$ , where  $\mathbf{y}_t \in \mathcal{Y} = \{0, 1\}^n$ . For any sequence of history  $(\mathbf{H}_t, \boldsymbol{\eta}_t)$  where  $\mathbf{H}_t = ([\mathbf{y}_1]_{I_1}, \dots, [\mathbf{y}_t]_{I_t}) = (h_1, \dots, h_t) \in \{0, 1\}^t$  and  $\boldsymbol{\eta}_t = (\pi_1, \dots, \pi_t) \in \{0, 1\}^t$ , let

$$\chi_{t,j}(\mathbf{H}_t, \boldsymbol{\eta}_t) = \mathbb{P}_j([\mathbf{y}_1]_{I_1} = h_1, \dots, [\mathbf{y}_t]_{I_t} = h_t, \boldsymbol{\eta}_t).$$

Note that  $\mathbb{P}_j$  is a measure on the  $\sigma$ -algebra  $\mathcal{A}$  as mentioned above, and the randomness comes from the adversary's moves (the adversary plays a randomized  $\mathbf{y}_t$  at time  $t$ , where the  $j^{\text{th}}$  coordinate of  $\mathbf{y}_t$  is a Bernoulli variable with mean either  $(1+\mu)/2$  or  $(1-\mu)/2$  depending on his choice of  $\zeta$  at the beginning of the game). From our assumption that the player is deterministic, for any  $\mathbf{H}_T \in \{0, 1\}^T$  and history of exploration  $\boldsymbol{\eta}_T$ . Then,

$$\mathbb{E}_i[T_j | [\mathbf{y}_1]_{I_1} = h_1, \dots, [\mathbf{y}_T]_{I_T} = h_T, \boldsymbol{\eta}_T] = \mathbb{E}_0[T_j | [\mathbf{y}_1]_{I_1} = h_1, \dots, [\mathbf{y}_T]_{I_T} = h_T, \boldsymbol{\eta}_T], \quad \forall 1 \leq i \leq n, \quad (11)$$

which means that no matter what the initial  $\zeta$  decided by the adversary is, the player has the same sequence of moves given the same history. Therefore,

$$\begin{aligned} \mathbb{E}_j[T_j] - \mathbb{E}_0[T_j] &= \sum_{\mathbf{H}_T \in \{0,1\}^T, \boldsymbol{\eta}_T \in \{0,1\}^T} \chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T) \mathbb{E}_j[T_j | [\mathbf{y}_1]_{I_1} = h_1, \dots, [\mathbf{y}_T]_{I_T} = h_T, \boldsymbol{\eta}_T] \\ &\quad - \sum_{\mathbf{H}_T \in \{0,1\}^T, \boldsymbol{\eta}_T \in \{0,1\}^T} \chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T) \mathbb{E}_0[T_j | [\mathbf{y}_1]_{I_1} = h_1, \dots, [\mathbf{y}_T]_{I_T} = h_T, \boldsymbol{\eta}_T] \\ &\stackrel{(1)}{=} \sum_{\mathbf{H}_T \in \{0,1\}^T, \boldsymbol{\eta}_T \in \{0,1\}^T} (\chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T) - \chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T)) \mathbb{E}_j[T_j | [\mathbf{y}_1]_{I_1} = h_1, \dots, [\mathbf{y}_T]_{I_T} = h_T, \boldsymbol{\eta}_T] \\ &\leq \sum_{(\mathbf{H}_T, \boldsymbol{\eta}_T) : \chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T) > \chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T)} (\chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T) - \chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T)) \mathbb{E}_j[T_j | [\mathbf{y}_1]_{I_1} = h_1, \dots, [\mathbf{y}_T]_{I_T} = h_T, \boldsymbol{\eta}_T] \\ &\stackrel{(2)}{\leq} T \sum_{(\mathbf{H}_T, \boldsymbol{\eta}_T) : \chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T) > \chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T)} (\chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T) - \chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T)), \end{aligned} \quad (12)$$

where Equation (1) follows from Equation (11) and Inequality (2) follows from  $\mathbb{E}_j[T_j | [\mathbf{y}_1]_{I_1} = h_1, \dots, [\mathbf{y}_T]_{I_T} = h_T, \boldsymbol{\eta}_T] \leq T$ . See that  $\sum_{j=1}^n \mathbb{E}_0[T_j] = T$  since on each round, player 1's action is in  $\{1, 2, \dots, n\}$ . We can bound the total variation using Pinsker's inequality:<sup>26</sup>

$$\sum_{(\mathbf{H}_T, \boldsymbol{\eta}_T) : \chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T) > \chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T)} (\chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T) - \chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T)) \leq \sqrt{\frac{1}{2} \text{KL}(\chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T) || \chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T))}. \quad (13)$$

Putting Equations (12) and (13) together and applying Jensen's inequality to the concave the square root function, we get

$$\begin{aligned} \frac{1}{n} \sum_{j=1}^n \mathbb{E}_j[T_j] &\geq \frac{1}{n} \sum_{j=1}^n \left( \mathbb{E}_0[T_j] + T \sum_{\mathbf{H}_T : \chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T) > \chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T)} (\chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T) - \chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T)) \right) \\ &\leq T \left( \frac{1}{n} + \frac{1}{n} \sum_{j=1}^n \sqrt{\frac{1}{2} \text{KL}(\chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T) || \chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T))} \right) \\ &\leq T \left( \frac{1}{n} + \sqrt{\frac{1}{2n} \sum_{j=1}^n \text{KL}(\chi_{T,0}(\mathbf{H}_T, \boldsymbol{\eta}_T) || \chi_{T,j}(\mathbf{H}_T, \boldsymbol{\eta}_T))} \right). \end{aligned} \quad (14)$$

Recall that from the definition of  $\chi_{t,j}$ , we have the following conditional distribution:

$$\chi_{t,j}(h_t | \mathbf{H}_{t-1}, \boldsymbol{\eta}_{t-1}) = \mathbb{P}_j([\mathbf{y}_t]_{I_t} = h_t | [\mathbf{y}_1]_{I_1} = h_1, \dots, [\mathbf{y}_{t-1}]_{I_{t-1}} = h_{t-1}, \boldsymbol{\eta}_{t-1}).$$

Applying the chain rule, we have

$$\begin{aligned} \text{KL}(\chi_{T,0} || \chi_{T,j}) &\stackrel{(1)}{=} \sum_{t=1}^T \sum_{\mathbf{H}_{t-1}, \boldsymbol{\eta}_{t-1}} \chi_{t-1,0}(\mathbf{H}_{t-1}, \boldsymbol{\eta}_{t-1}) \text{KL}(\chi_{t,0}(\cdot | \mathbf{H}_{t-1}, \boldsymbol{\eta}_{t-1}) || \chi_{t,j}(\cdot | \mathbf{H}_{t-1}, \boldsymbol{\eta}_{t-1})) \\ &\stackrel{(2)}{=} \sum_{t=1}^T \sum_{\mathbf{H}_{t-1}, \boldsymbol{\eta}_{t-1}} \chi_{t-1,0}(\mathbf{H}_{t-1}, \boldsymbol{\eta}_{t-1}) \mathbb{1}(\{I_t = j \text{ and } \pi_t = 1 | \mathbf{H}_{t-1}, \boldsymbol{\eta}_{t-1}\}) \text{KL}\left(\frac{1-\mu}{2} \middle| \middle| \frac{1+\mu}{2}\right) \\ &\stackrel{(3)}{=} \text{KL}\left(\frac{1-\mu}{2} \middle| \middle| \frac{1+\mu}{2}\right) \mathbb{E}_0 \left[ \sum_{t=1}^T \mathbb{1}(\{I_t = j \text{ and } \pi_t = 1\}) \right]. \end{aligned}$$

Here, Equation (1) follows from applying the chain rule to  $\chi_{T,0}$  and  $\chi_{T,j}$ . Equation (2) holds because  $\chi_{t,0}(\cdot | \mathbf{H}_{t-1}, \boldsymbol{\eta}_{t-1}) = \text{Ber}\left(\frac{1-\mu}{2}\right)$  and  $\chi_{t,j}(\cdot | \mathbf{H}_{t-1}, \boldsymbol{\eta}_{t-1}) = \text{Ber}\left(\frac{1+\mu}{2}\right)$  when we play the arm  $j$  on round  $t$ ,  $I_t = j$ , and observe the payoff,  $\pi_t = 1$ . Otherwise, they are identical and we have  $\text{KL}(\chi_{t,0}(\cdot | \mathbf{H}_{t-1}, \boldsymbol{\eta}_{t-1}) || \chi_{t,j}(\cdot | \mathbf{H}_{t-1}, \boldsymbol{\eta}_{t-1})) = 0$ . Lastly, we get Equation (3) by factoring out  $\text{KL}\left(\frac{1-\mu}{2} \middle| \middle| \frac{1+\mu}{2}\right)$ , and collecting all the probability terms ( $\chi_{t,0}(\mathbf{H}_t, \boldsymbol{\mu}_t)$  for all  $t$ ) to form the expectation of  $\mathbb{1}(\{I_t = j \text{ and } \pi_t = 1\})$  with respect to  $\mathbb{P}_0$ .

<sup>26</sup> Pinsker's inequality bounds the total variation distance in terms of KL divergence. For two probability distributions  $P$  and  $Q$ , Pinsker's inequality states that  $\|P - Q\|_{TV} \leq \sqrt{\frac{1}{2} \text{KL}(P||Q)}$  where  $\|P - Q\|_{TV}$  is the total variation distance  $\sup_A \{|P(A) - Q(A)|\}$  over measurable events  $A$ . Taking  $A = \{x : P(x) > Q(x)\}$ , we get  $\sum_{x: P(x) > Q(x)} |P(x) - Q(x)| \leq \sqrt{\frac{1}{2} \text{KL}(P||Q)}$ . See Section A.2 Cesa-Bianchi and Lugosi (2006) for details.

Summing over  $j$  and applying  $\text{KL}(p||q) \leq \frac{(p-q)^2}{q(1-q)}$ :

$$\begin{aligned} \sum_{j=1}^n \text{KL}(\chi_{T,0}||\chi_{T,j}) &= \text{KL}\left(\frac{1-\mu}{2} \middle| \middle| \frac{1+\mu}{2}\right) \sum_{j=1}^n \mathbb{E}_0 \left[ \sum_{t=1}^T \mathbb{1}(\{I_t = j \text{ and } \pi_t = 1\}) \right] \\ &= \text{KL}\left(\frac{1-\mu}{2} \middle| \middle| \frac{1+\mu}{2}\right) \mathbb{E}_0 \left[ \sum_{t=1}^T \sum_{j=1}^n \mathbb{1}(\{I_t = j \text{ and } \pi_t = 1\}) \right] \\ &= \text{KL}\left(\frac{1-\mu}{2} \middle| \middle| \frac{1+\mu}{2}\right) \mathbb{E}_0 \left[ \sum_{t=1}^T \mathbb{1}(\{\pi_t = 1\}) \right] \leq \frac{4\mu^2}{1-\mu^2} M, \end{aligned} \quad (15)$$

where the last line follows from the assumption that the number of rounds the player explores is  $M$ .

Putting Equation (10), (14) and (15) altogether we get:

$$\mathbb{E}_j \left[ d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) \right] \geq \mu \left( 1 - \frac{1}{n} - \sqrt{\frac{1}{2n} \frac{4\mu^2}{1-\mu^2} M} \right) \geq \mu \left( 1 - \frac{1}{n} - 4\mu \sqrt{\frac{M}{6n}} \right), \quad (16)$$

where the last inequality follows from  $\mu \leq 1/4$ . Taking  $\mu = \lambda \sqrt{\frac{n}{M}}$ , we have:

$$\mathbb{E}_j \left[ d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) \right] \geq \lambda \sqrt{\frac{n}{M}} \left( \frac{1}{2} - \frac{4\lambda}{\sqrt{6}} \right) \geq \Omega \left( \frac{1}{\sqrt{M}} \right) = \Omega \left( \frac{D}{\sqrt{M}} \right), \quad (17)$$

where the last equality follows from  $\underline{D} \leq D(\mathbf{p}) = 1$  in this case since the adversary's move is in  $\{0, 1\}^n$ . We finish the proof by choosing the constant  $\lambda$  to be small enough to ensure that  $\left(\frac{1}{2} - \frac{4\lambda}{\sqrt{6}}\right)$  is positive. ■

## Appendix E: Application to Non-monotone (Continuous) Submodular Maximization

*Problem definition.* Consider the NON-MONOTONE SUBMODULAR MAXIMIZATION (NSM) problem, for both set and continuous functions. For set functions, our goal is to maximize a non-monotone submodular set function without any constraints, and for continuous functions, our goal is to maximize a non-monotone continuous submodular function, either weak-DR or strong-DR, over the unit hypercube  $[0, 1]^n$ ; see the definition of weak-DR and strong-DR continuous submodular functions below.

*Continuous submodular functions.* We defined set submodular functions in Example 1. The concept of submodularity can be extended from subset lattice (above definition) to any discrete or continuous lattice. In particular, by considering the positive orthant cone lattice, we can define the continuous variant of set submodularity. A continuous multivariate function  $f: [0, 1]^n \rightarrow [0, 1]$  is submodular if for all  $\mathbf{x}, \mathbf{y} \in [0, 1]^n$ ,

$$f(\mathbf{x} \vee \mathbf{y}) + f(\mathbf{x} \wedge \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y}),$$

where  $\vee$  and  $\wedge$  are coordinate-wise max and min operations. As an equivalent definition (Bian et al. 2017),  $f$  is continuous submodular if for all  $i \in [n]$ ,  $z \in [0, 1]$ ,  $\mathbf{x}_{-i} \preceq \mathbf{y}_{-i} \in [0, 1]^{n-1}$ , and  $\delta \geq 0$ ,

$$f(z + \delta, \mathbf{x}_{-i}) - f(z, \mathbf{x}_{-i}) \geq f(z + \delta, \mathbf{y}_{-i}) - f(z, \mathbf{y}_{-i}).$$

The above class of continuous functions is also referred to as the weak-Diminishing Return (weak-DR) submodular in the literature Wolsey (1982), Bach (2019). We further consider a special subclass of these

functions satisfying concavity along each coordinate, referred to as the strong-Diminishing Return (strong-DR). A continuous multivariate function  $f : [0, 1]^n \rightarrow [0, 1]$  is strong-DR submodular if for all  $i \in [n]$ ,  $\mathbf{x} \preceq \mathbf{y} \in [0, 1]^n$ , and  $\delta \geq 0$ ,

$$f(x_i + \delta, \mathbf{x}_{-i}) - f(\mathbf{x}) \geq f(y_i + \delta, \mathbf{y}_{-i}) - f(\mathbf{y}),$$

where  $\mathbf{x}_{-i}$  (resp.  $\mathbf{y}_{-i}$ ) is an  $(n-1)$ -dimensional vector with all coordinate values of  $\mathbf{x}$  (resp.  $\mathbf{y}$ ) except  $i$ , and  $\mathbf{x} \preceq \mathbf{y}$  if and only if  $\forall j \in [n] : x_j \leq y_j$ .

*Offline problems for submodular functions.* For set functions, the offline algorithm of Buchbinder et al. 2015 gives a  $1/2$ -approximation factor, which is known to be the best possible approximation factor with polynomial query calls to the function (Feige et al. 2011). For the continuous case, under both weak-DR and strong-DR submodularity, the offline algorithm of Niazadeh et al. 2018 gives a  $1/2$ -approximation factor for Lipschitz continuous functions, which again achieves the best possible approximation factor with polynomial query calls to the function.

To have a unified offline problem and algorithm capturing both of the above variations, we first consider a slight reformulation where a continuous (weak-DR) submodular function is restricted to a discrete domain  $\mathcal{R}^n$  instead of  $[0, 1]^n$ . Here,  $\mathcal{R} = \{\rho_1, \rho_2, \dots, \rho_m\}$  is the finite set of possible coordinate values, where  $|\mathcal{R}| = m$  and  $\rho_1 < \rho_2 < \dots < \rho_m$  are real numbers. Note that  $\mathcal{R} = \{0, 1\}$  when we focus on set functions. For Lipschitz continuous functions, one should think of  $\mathcal{R}^n$  as an  $\epsilon$ -net that discretizes the function with  $O(\epsilon)$  additive error due to Lipschitzness.

Given this unified setting, we essentially consider discrete functions  $f : \mathcal{R}^n \rightarrow [0, 1]$  that satisfy a discrete version of (weak-DR) submodularity. This property is exactly the same as continuous submodularity, with a slight modification that we only consider points  $\mathbf{x} \in \mathcal{R}^n$ . Given such a function, our goal in the offline problem is to solve the optimization problem  $\max_{\mathbf{z} \in \mathcal{R}^n} f(\mathbf{z})$ . We refer to this problem as *discretized submodular maximization*. Note that this problem is an instance of problem (1), where both  $\mathcal{D}$  and the feasible region  $\mathcal{C}$  are  $\mathcal{R}^n$ , and our function class is the class of submodular functions  $f$  described above.

Inspired by the algorithms in Buchbinder et al. (2015) and Niazadeh et al. (2018), we then present a unified offline algorithm (which essentially is an adaptation of the algorithm in Niazadeh et al. (2018) restricted to the discrete domain  $\mathcal{R}^n$ ) with the same  $1/2$ -approximation factor for the proposed unified offline problem. This is presented in Algorithm 7. We then transform this offline algorithm to online full-information and bandit learning algorithms using our framework.

Algorithm 7 is a modified version of the continuous randomized bi-greedy algorithm by Niazadeh et al. (2018). The difference between Algorithm 7 and the continuous randomized bi-greedy algorithm is discussed in Section I.1 in the appendix. Throughout this section, we use the notation  $(z', \mathbf{z}_{-i})$  to denote the point constructed by taking  $\mathbf{z}$  and replacing its  $i^{\text{th}}$  coordinate value with  $z'$ , and  $f(z', \mathbf{z}_{-i})$  to denote the function evaluated at the corresponding point. The algorithm keeps track of two points: *lower bound*  $\underline{\mathbf{z}}^{(i)}$  and *upper bound*  $\bar{\mathbf{z}}^{(i)}$ , where initially,  $\underline{\mathbf{z}}^{(0)} = (\rho_1, \dots, \rho_1)$ , and  $\bar{\mathbf{z}}^{(0)} = (\rho_m, \dots, \rho_m)$ . The lower and upper bounds get updated as the algorithm goes through  $n$  subproblems. In subproblem  $i$ , the algorithm decides about the  $i^{\text{th}}$  coordinate: it sets the  $i^{\text{th}}$  coordinate to  $z'_i$ , where  $z'_i$  is drawn from distribution  $\theta^{(i)} \in \Delta(\mathcal{R})$ . Here, this distribution is chosen in a way to satisfy the following condition  $\mathbb{E}_{z' \sim \theta^{(i)}} \left[ \frac{1}{2} \alpha^{(i)}(z') + \frac{1}{2} \beta^{(i)}(z') - \zeta^{(i)}(\hat{\mathbf{z}}, z') \right] \geq 0$

for all  $\hat{z} \in \mathcal{R}$ . Note that  $\alpha^{(i)}(z') = f(z', \underline{z}_{-i}^{(i-1)}) - f(\underline{z}^{(i-1)})$  is the marginal value of increasing the value of  $i^{\text{th}}$ -coordinate from  $\rho_1$  to  $z'$  when the rest of coordinates are  $\underline{z}_{-i}^{(i-1)}$ , and similarly  $\beta^{(i)}(z') = f(z', \bar{z}_{-i}^{(i-1)}) - f(\bar{z}^{(i-1)})$  is the marginal value of decreasing the  $i^{\text{th}}$  coordinate from  $\rho_m$  to  $z'$  when the rest of coordinates are  $\bar{z}_{-i}^{(i-1)}$ . Moreover,  $\zeta^{(i)}(\hat{z}, z')$  is equal to  $\alpha^{(i)}(\hat{z}) - \alpha^{(i)}(z')$  if  $\hat{z} \geq z'$  and  $\beta^{(i)}(\hat{z}) - \beta^{(i)}(z')$  otherwise. Roughly speaking,  $\zeta^{(i)}(\hat{z}, z')$  measures the extent to which setting the  $i^{\text{th}}$  coordinate to  $z'$ , rather than  $\hat{z}$ , is locally suboptimal. With this interpretation, the aforementioned condition ensures that the algorithm's choice for the  $i^{\text{th}}$  coordinate approximately compensates for the cost caused by the suboptimality of this choice. We refer the readers to Niazadeh et al. (2018) for a more detailed discussion on the intuition behind this condition.

We now show how to cast the above algorithm as an instance of OFFLINE-IG (Algorithm 1). In the language of Algorithm 1, the aforementioned condition can be presented using the following PAYOFF function:

$$j \in [m]: \quad [\text{PAYOFF}(\theta^{(i)}, \underline{z}^{(i-1)}, f)]_j = \mathbb{E}_{z' \sim \theta^{(i)}} \left[ \frac{1}{2} \alpha^{(i)}(z') + \frac{1}{2} \beta^{(i)}(z') - \zeta^{(i)}(\rho_j, z') \right] \geq 0. \quad (18)$$

Moreover, we have  $\Theta = \Delta(\mathcal{R})$ ,  $d_{\text{param}} = |\mathcal{R}| = m$ , and  $\mathbf{z}$  is the vector  $\underline{z}$  that starts as  $(\rho_1, \dots, \rho_1)^T$  then gets updated at each iteration.<sup>27</sup>

**THEOREM 8 (Online learning for discretized non-monotone submodular maximization).** *Let  $n$  be the number of dimensions and  $\mathcal{R} = \{\rho_1, \rho_2, \dots, \rho_m\}$  be the set of potential values that each coordinate  $i \in [n]$  can take. Assume that the maximum function value is normalized to one. Then, for the problem of maximizing a non-monotone submodular function in the online full-information setting, there exists a learning algorithm that obtains  $O(nT^{1/2}(\log m)^{1/2})$   $\frac{1}{2}$ -regret, where  $T$  is the number of rounds. Furthermore, in the online bandit setting, there exists an online learning algorithm that obtains  $O(nm^{2/3}T^{2/3}(\log m)^{1/3})$   $\frac{1}{2}$ -regret. Here, in both online algorithms, the benchmark in the regret bounds is  $\frac{1}{2} \max_{\mathbf{z} \in \mathcal{R}^n} \sum_{t=1}^T f_t(\mathbf{z})$ .*

The proof of Theorem 8, which is presented in Section E.1, has two main steps. In the first step, we show that the offline Algorithm 7 is an extended  $(\frac{1}{2}, \frac{1}{2})$ -robust approximation algorithm and in the second step, we show that it is bandit Blackwell reducible. The challenging part of the proof is to construct an explore sampling device that leads to an unbiased estimator for the payoff function. We then invoke Theorems 2 and 4 to get the final regret bounds.

The following is an immediate corollary of Theorem 8.

**COROLLARY 3 (Online learning for non-monotone set submodular maximization).** *Let  $n$  be the number of items, and assume the maximum function value is normalized to one. Then for the problem of maximizing a nonmonotone (set) submodular function in the online full-information setting, there exists an online learning algorithm that obtains  $O(n\sqrt{T})$   $\frac{1}{2}$ -regret, where  $T$  is the number of rounds. Furthermore, in the online bandit setting, there exists a learning algorithm that obtains  $O(nT^{2/3})$   $\frac{1}{2}$ -regret, where  $T$  is the number of rounds.*

<sup>27</sup> For any  $i \in [n]$ , one can construct  $\bar{z}^{(i)}$  from  $\underline{z}^{(i)}$  by replacing its last  $n-1$  coordinates with  $\rho_m$ . Thus, it suffices to define PAYOFF as a function of  $\underline{z}^{(i)}$ .

**Algorithm 7:** Greedy Algorithm for Discretized NSM (Niazadeh et al. 2018)**Input:** Discrete submodular function  $f$ .**Output:** Point  $\mathbf{z} \in \mathcal{R}^n$ .Set initial *lower bound*  $\underline{\mathbf{z}}^{(0)} \leftarrow (\rho_1, \rho_1, \dots, \rho_1)^T$  and *upper bound*  $\bar{\mathbf{z}}^{(0)} \leftarrow (\rho_m, \rho_m, \dots, \rho_m)^T$ .**for** coordinate  $i = 1, 2, \dots, n$  **do**Define the lower marginal function  $\alpha^{(i)} : \mathcal{R} \rightarrow [-1, +1]$  as

$$\alpha^{(i)}(z') = f(z', \underline{\mathbf{z}}_{-i}^{(i-1)}) - f(\underline{\mathbf{z}}^{(i-1)}).$$

Define the upper marginal function  $\beta^{(i)} : \mathcal{R} \rightarrow [-1, +1]$  as

$$\beta^{(i)}(z') = f(z', \bar{\mathbf{z}}_{-1}^{(i-1)}) - f(\bar{\mathbf{z}}^{(i-1)}).$$

Define comparison function  $\zeta^{(i)} : \mathcal{R} \times \mathcal{R} \rightarrow [-1, +1]$  as

$$\zeta^{(i)}(\hat{z}, z') = \begin{cases} \alpha^{(i)}(\hat{z}) - \alpha^{(i)}(z') & \text{if } \hat{z} \geq z' \\ \beta^{(i)}(\hat{z}) - \beta^{(i)}(z') & \text{if } \hat{z} \leq z' \end{cases}.$$

Local Optimization StepChoose  $\boldsymbol{\theta}^{(i)} \in \Delta(\mathcal{R})$  so that for all  $\hat{z} \in \mathcal{R}$ ,

$$\mathbb{E}_{z' \sim \boldsymbol{\theta}^{(i)}} \left[ \frac{1}{2} \alpha^{(i)}(z') + \frac{1}{2} \beta^{(i)}(z') - \zeta^{(i)}(\hat{z}, z') \right] \geq 0 \quad (19)$$

(done in Niazadeh et al. (2018) via preprocessing and computing a 2D convex hull).

Local Update StepSample  $z'_i \sim \boldsymbol{\theta}^{(i)}$ . Set  $\underline{\mathbf{z}}^{(i)} \leftarrow \underline{\mathbf{z}}^{(i-1)}$  and  $\bar{\mathbf{z}}^{(i)} \leftarrow \bar{\mathbf{z}}^{(i-1)}$  and then update their  $i^{\text{th}}$  coordinate:

$$[\underline{\mathbf{z}}^{(i)}]_i \leftarrow z'_i \text{ and } [\bar{\mathbf{z}}^{(i)}]_i \leftarrow z'_i.$$

**return**  $\mathbf{z} \leftarrow \underline{\mathbf{z}}^{(n)}$ .

So far, we assumed that for the continuous submodular functions, the set of potential value for each coordinate is finite and belongs to set  $\mathcal{R} = \{\rho_1, \rho_2, \dots, \rho_m\}$ , rather than the interval  $[0, 1]$ , and we design learning algorithms with sublinear regret bounds where the regrets are computed with respect to  $\frac{1}{2} \max_{\mathbf{z} \in \mathcal{R}} \sum_{t=1}^T f_t(\mathbf{z})$ . Now, one may wonder if one can design learning algorithms against the benchmark of  $\frac{1}{2} \max_{\mathbf{z} \in [0, 1]^n} \sum_{t=1}^T f_t(\mathbf{z})$  that allows the coordinates to be any number in  $[0, 1]^n$ . The following corollary answers this question for any  $L$ -Lipschitz non-monotone continuous submodular functions.

**COROLLARY 4 (Online learning for  $L$ -Lipschitz continuous submodular maximization).** *Let  $n$  be the number of dimensions, and assume the maximum function value is normalized to one. Then for the problem of maximizing a coordinate-wise  $L$ -Lipschitz non-monotone (continuous) submodular function in the online full-information setting, there exists a learning algorithm that obtains  $O(nT^{1/2} \log^{1/2}(LT))$   $\frac{1}{2}$ -regret, where  $T$  is the number of rounds. Furthermore, in the online bandit setting, there exists a learning algorithm that obtains  $O(nL^{2/5}T^{4/5} \log^{1/3}(LT))$   $\frac{1}{2}$ -regret, where  $T$  is the number of rounds. Here, in both online algorithms, the benchmark in the regret bounds is  $\frac{1}{2} \max_{\mathbf{z} \in [0, 1]^n} \sum_{t=1}^T f_t(\mathbf{z})$ .*

Proofs of the above corollaries are in Section E.2.

### E.1. Proof of Theorem 8

*Proof.* We will show that our meta Algorithms 2 and 4 work by verifying the following conditions.

(i) *Algorithm 6 is an extended  $(\frac{1}{2}, \frac{1}{2})$ -robust approximation algorithm.* Following the analysis of the bi-greedy algorithm in Buchbinder et al. (2015), we consider three sequences of points: the lower bound sequence  $\underline{z}^{(i)}$ , the upper bound sequence  $\bar{z}^{(i)}$ , and the hybrid-optimal sequence  $\mathbf{z}^{*(i)}$ . The key proof idea is to bound the decrease in the hybrid-optimal sequence value  $\mathbf{z}^{*(i)}$  with the total increase in the lower bound and upper bound sequence values. We define  $\underline{z}^{(i)}$  and  $\bar{z}^{(i)}$  to agree on the first  $i$  coordinates, while the rest of the coordinates are  $\rho_1$  for  $\underline{z}^{(i)}$  and  $\rho_m$  for  $\bar{z}^{(i)}$ . The hybrid-optimal sequence starts from  $\mathbf{z}^{*(0)} \triangleq \mathbf{z}^*$ , then  $\mathbf{z}^{*(i)}$  is equal to  $\mathbf{z}^{*(i-1)}$  but with the  $i^{\text{th}}$  coordinate replaced with the sampled  $z'_i \sim \theta^{(i)}$ .

Importantly, if the  $i^{\text{th}}$ -coordinate of the optimal vector  $\mathbf{z}^*$ , which is  $z_i^*$ , is less than our sampled point  $z'_i$  from the  $i^{\text{th}}$  subproblem/iteration, then the loss in value of the hybrid-optimal sequence is bounded by a difference of two  $\beta^{(i)}$  evaluations. In particular, the submodularity of  $f$  implies:

$$\begin{aligned} f(z_i^*, \mathbf{z}_{-i}^{*(i-1)}) + f(z'_i, \bar{\mathbf{z}}_{-i}^{(i-1)}) &\leq f(z'_i, \mathbf{z}_{-i}^{*(i-1)}) + f(z_i^*, \bar{\mathbf{z}}_{-i}^{(i-1)}) \\ f(\mathbf{z}^{*(i-1)}) + \beta^{(i)}(z'_i) &\leq f(\mathbf{z}^{*(i)}) + \beta^{(i)}(z_i^*) \\ f(\mathbf{z}^{*(i-1)}) - f(\mathbf{z}^{*(i)}) &\leq \beta^{(i)}(z_i^*) - \beta^{(i)}(z'_i). \end{aligned}$$

There is also the symmetric case where the  $i^{\text{th}}$ -coordinate of the optimal vector  $z_i^*$  is greater than our sampled point  $z'_i$  from the  $i^{\text{th}}$  subproblem:

$$\begin{aligned} f(z'_i, \underline{\mathbf{z}}_{-i}^{(i-1)}) + f(z_i^*, \mathbf{z}_{-i}^{*(i-1)}) &\leq f(z_i^*, \underline{\mathbf{z}}_{-i}^{(i-1)}) + f(z'_i, \mathbf{z}_{-i}^{*(i-1)}) \\ \alpha^{(i)}(z'_i) + f(\mathbf{z}^{*(i-1)}) &\leq \alpha^{(i)}(z_i^*) + f(\mathbf{z}^{*(i)}) \\ f(\mathbf{z}^{*(i-1)}) - f(\mathbf{z}^{*(i)}) &\leq \alpha^{(i)}(z_i^*) - \alpha^{(i)}(z'_i). \end{aligned}$$

Combining the two cases yields (this inequality explains our definition of  $\zeta^{(i)}$ ):

$$f(\mathbf{z}^{*(i-1)}) - f(\mathbf{z}^{*(i)}) \leq \zeta^{(i)}(z_i^*, z'_i). \quad (20)$$

Also, just by the definition of  $\alpha^{(i)}$  and  $\beta^{(i)}$  we know that:

$$f(\underline{\mathbf{z}}^{(i)}) - f(\underline{\mathbf{z}}^{(i-1)}) = \alpha^{(i)}(z'_i) \quad (21)$$

$$f(\bar{\mathbf{z}}^{(i)}) - f(\bar{\mathbf{z}}^{(i-1)}) = \beta^{(i)}(z'_i). \quad (22)$$

We are now ready to consider the  $\theta_t^{(i)}$  which guarantees that for all  $i \in [n]$  and  $\hat{z} \in \mathcal{R}$ , including  $z_i^*$ :

$$\sum_{t=1}^T \mathbb{E}_{z'_i \sim \theta_t^{(i)}} \left[ \frac{1}{2} \alpha_t^{(i)}(z'_i) + \frac{1}{2} \beta_t^{(i)}(z'_i) - \zeta_t^{(i)}(z_i^*, z'_i) \right] \geq -h(T).$$

Note that  $\alpha_t^{(i)}$ ,  $\beta_t^{(i)}$ , and  $\zeta_t^{(i)}$  are respectively obtained by replacing  $f$  with  $f_t$  in the definition of  $\alpha^{(i)}$ ,  $\beta^{(i)}$ , and  $\zeta^{(i)}$ . We sum those inequalities together and then apply Equations (20), (21), and (22):

$$-nh(T)$$

$$\begin{aligned}
&\leq \sum_{t=1}^T \sum_{i=1}^n \mathbb{E}_{z'_t \sim \theta_t^{(i)}} \left[ \frac{1}{2} \alpha_t^{(i)}(z'_t) + \frac{1}{2} \beta_t^{(i)}(z'_t) - \zeta_t^{(i)}(z'_t, z'_t) \right] \\
&\leq \sum_{t=1}^T \sum_{i=1}^n \mathbb{E} \left[ \frac{1}{2} [f_t(\mathbf{z}^{(i)}) - f_t(\mathbf{z}^{(i-1)})] + \frac{1}{2} [f_t(\bar{\mathbf{z}}^{(i)}) - f_t(\bar{\mathbf{z}}^{(i-1)})] - [f_t(\mathbf{z}^{*(i-1)}) - f_t(\mathbf{z}^{*(i)})] \right] \\
&= \sum_{t=1}^T \mathbb{E} \left[ \frac{1}{2} [f_t(\mathbf{z}^{(n)}) - f_t(\mathbf{z}^{(0)})] + \frac{1}{2} [f_t(\bar{\mathbf{z}}^{(n)}) - f_t(\bar{\mathbf{z}}^{(0)})] - [f_t(\mathbf{z}^{*(0)}) - f_t(\mathbf{z}^{*(n)})] \right] \\
&= \sum_{t=1}^T \mathbb{E} \left[ \frac{1}{2} \left[ f_t(\mathbf{z}_t) - \underbrace{f_t(\mathbf{z}^{(0)})}_{\geq 0} \right] + \frac{1}{2} \left[ f_t(\mathbf{z}_t) - \underbrace{f_t(\bar{\mathbf{z}}^{(0)})}_{\geq 0} \right] - [f_t(\mathbf{z}^*) - f_t(\mathbf{z}_t)] \right] \\
&\leq \sum_{t=1}^T \mathbb{E} [2f_t(\mathbf{z}_t) - f_t(\mathbf{z}^*)].
\end{aligned}$$

See that the fourth equality is because the algorithm returns  $\mathbf{z}_t = \mathbf{z}^{(n)} = \bar{\mathbf{z}}^{(n)}$  at round  $t$ . We finish by moving terms between sides and dividing by two:

$$\begin{aligned}
\sum_{t=1}^T \mathbb{E} [2f_t(\mathbf{z}_t) - f_t(\mathbf{z}^*)] &\geq -nh(T) \\
\sum_{t=1}^T \mathbb{E} [f_t(\mathbf{z}_t)] &\geq \frac{1}{2} \sum_{t=1}^T f_t(\mathbf{z}^*) - \frac{1}{2} nh(T).
\end{aligned}$$

Thus, our algorithm is an extended  $(\frac{1}{2}, \frac{1}{2})$ -robust approximation.

(ii) *Algorithm 7 is bandit Blackwell reducible.* We first show that Algorithm 7 is Blackwell reducible. Consider an instance  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  of Blackwell where  $\mathcal{X} \triangleq \Theta = \Delta(\mathcal{R})$  and  $\mathcal{Y} \triangleq \Delta(\mathcal{C} \times \mathcal{F}) = \Delta(\mathcal{R}^{[n]} \times \mathcal{F})$ . Our synthetic Blackwell adversary function is the deterministic distribution that has weight 1 on its input (point, function) pair and 0 anywhere else, i.e.  $\text{AdvB}(\mathbf{z}, f) = \kappa$  where  $\kappa(\mathbf{z}, f) = 1$ . The (asymmetric) biaffine Blackwell payoff  $\mathbf{p}$  is the expectation of the PAYOFF function from Equation (18) over its second input:

$$\mathbf{p}(\boldsymbol{\theta}, \kappa) \triangleq \mathbb{E}_{(\mathbf{z}, f) \sim \kappa} [\text{PAYOFF}(\boldsymbol{\theta}, \mathbf{z}, f)].$$

The positive orthant  $S$  is response-satisfiable since given any player 2 distribution  $\kappa$  over (point, function) pairs, we can convert each pair into the marginal functions  $\alpha^{(i)}$  and  $\beta^{(i)}$ . Averaging these marginal functions together according to their likelihood in  $\kappa$  does not impact the submodularity fact we require for our proofs. We can think of  $\mathbf{p}(\boldsymbol{\theta}, \kappa)$  as

$$\mathbf{p}(\boldsymbol{\theta}, \kappa) \triangleq \mathbb{E}_{(\mathbf{z}, f) \sim \kappa} [\text{PAYOFF}(\boldsymbol{\theta}, \mathbf{z}, f)] = \text{PAYOFF}(\boldsymbol{\theta}, \mathbf{z}, f'),$$

for another submodular function  $f' \in \mathcal{F}$  because a weighted average of submodular function is submodular. Since for any submodular functions  $f \in \mathcal{F}$  and  $\mathbf{z} \in \mathcal{C}$ , we show that we can find  $\boldsymbol{\theta}$  such that  $\text{PAYOFF}(\boldsymbol{\theta}, \mathbf{z}, f) \geq 0$ , for any  $\kappa$ , the algorithm can find  $\boldsymbol{\theta}$  such that  $\mathbf{p}(\boldsymbol{\theta}, \kappa)$  is in  $S$ . Therefore, Algorithm 7 is Blackwell reducible.

To show that Algorithm 7 is bandit Blackwell reducible, we need to construct an unbiased estimator for  $\mathbf{p}$  and an explore sampling device  $U$ . In subproblem  $i$ ,  $U$  receives pairs of the form  $(\boldsymbol{\theta}, \mathbf{z}^{(i-1)})$  and returns  $(\mathbf{w}_{\text{exp}}, \mathbf{z}_{\text{exp}})$  such that (i) for all  $f \in \mathcal{F}, \boldsymbol{\theta} \in \Theta, \mathbf{z}^{(i-1)} \in \mathcal{D}, \hat{\mathbf{p}}(\boldsymbol{\theta}, \text{AdvB}(\mathbf{z}^{(i-1)}, f)) = f(\mathbf{z}_{\text{exp}}) \mathbf{w}_{\text{exp}}$  where  $(\mathbf{w}_{\text{exp}}, \mathbf{z}_{\text{exp}}) \sim U(\boldsymbol{\theta}, \mathbf{z}^{(i-1)})$ , and (ii)  $\hat{\mathbf{p}}$  is an unbiased estimator for the actual payoff, i.e.  $\forall \boldsymbol{\theta} \in \Theta, \kappa \in \mathcal{Y}$ , we have  $\mathbb{E}[\hat{\mathbf{p}}(\boldsymbol{\theta}, \kappa)] = \mathbf{p}(\boldsymbol{\theta}, \kappa)$ .

Because we would like to construct an unbiased estimator of the actual payoff  $\mathbf{p}$ , which is an expectation (over  $\kappa$ ) of the payoff function PAYOFF, which is further an affine combination of the functions  $\alpha^{(i)}, \beta^{(i)}$ , and  $\zeta^{(i)}$  on  $\mathcal{R}$ , we construct unbiased estimators from function evaluations for these functions. Observe that given  $\mathbf{z}^{(i-1)}$ ,  $U$  can immediately reconstruct the corresponding upper bound point:

$$\bar{\mathbf{z}}^{(i-1)} \leftarrow \mathbf{z}^{(i-1)} \vee \left( \underbrace{\rho_1, \dots, \rho_1}_{\text{first } (i-1) \text{ coordinates}}, \underbrace{\rho_m, \dots, \rho_m}_{\text{last } (n-i+1) \text{ coordinates}} \right)^T = \left( z'_1, \dots, z'_{i-1}, \underbrace{\rho_m, \dots, \rho_m}_{\text{last } (n-i+1) \text{ coordinates}} \right)^T.$$

We can use  $\mathbf{z}^{(i-1)}$  and  $\bar{\mathbf{z}}^{(i-1)}$  to express the marginal functions  $\alpha^{(i)}$  and  $\beta^{(i)}$ ,

$$\alpha^{(i)} \triangleq \begin{bmatrix} \alpha^{(i)}(\rho_1) \\ \alpha^{(i)}(\rho_2) \\ \vdots \\ \alpha^{(i)}(\rho_m) \end{bmatrix} = \begin{bmatrix} f(\rho_1, \bar{\mathbf{z}}_{-i}^{(i-1)}) - f(\mathbf{z}^{(i-1)}) \\ f(\rho_2, \bar{\mathbf{z}}_{-i}^{(i-1)}) - f(\mathbf{z}^{(i-1)}) \\ \vdots \\ f(\rho_m, \bar{\mathbf{z}}_{-i}^{(i-1)}) - f(\mathbf{z}^{(i-1)}) \end{bmatrix}$$

$$\beta^{(i)} \triangleq \begin{bmatrix} \beta^{(i)}(\rho_1) \\ \beta^{(i)}(\rho_2) \\ \vdots \\ \beta^{(i)}(\rho_m) \end{bmatrix} = \begin{bmatrix} f(\rho_1, \bar{\mathbf{z}}_{-i}^{(i-1)}) - f(\bar{\mathbf{z}}^{(i-1)}) \\ f(\rho_2, \bar{\mathbf{z}}_{-i}^{(i-1)}) - f(\bar{\mathbf{z}}^{(i-1)}) \\ \vdots \\ f(\rho_m, \bar{\mathbf{z}}_{-i}^{(i-1)}) - f(\bar{\mathbf{z}}^{(i-1)}) \end{bmatrix}.$$

These can be used in turn to express our comparison function  $\zeta^{(i)}$ :

$$\zeta^{(i)} \triangleq \begin{bmatrix} \zeta^{(i)}(\rho_1, \rho_1) & \zeta^{(i)}(\rho_1, \rho_2) & \cdots & \zeta^{(i)}(\rho_1, \rho_m) \\ \zeta^{(i)}(\rho_2, \rho_1) & \zeta^{(i)}(\rho_2, \rho_2) & \cdots & \zeta^{(i)}(\rho_2, \rho_m) \\ \vdots & \vdots & \ddots & \vdots \\ \zeta^{(i)}(\rho_m, \rho_1) & \zeta^{(i)}(\rho_m, \rho_2) & \cdots & \zeta^{(i)}(\rho_m, \rho_m) \end{bmatrix}$$

$$= \text{diag}(\alpha^{(i)}) \mathbf{L}_{m,m} - \mathbf{L}_{m,m} \text{diag}(\alpha^{(i)}) + \text{diag}(\beta^{(i)}) \mathbf{U}_{m,m} - \mathbf{U}_{m,m} \text{diag}(\beta^{(i)}),$$

where  $\mathbf{L}_{m,m}$  is the lower-triangular matrix defined by  $[\mathbf{L}_{m,m}]_{i,j} = \mathbb{1}[i > j]$  and  $\mathbf{U}_{m,m}$  is the upper-triangular matrix defined by  $[\mathbf{U}_{m,m}]_{i,j} = \mathbb{1}[i < j]$ . Our desired payoff function can be expressed using all three of these functions:

$$\text{PAYOFF}(\boldsymbol{\theta}, \mathbf{z}^{(i-1)}, f) = \left[ \frac{1}{2} \mathbf{1}_m (\alpha^{(i)})^T + \frac{1}{2} \mathbf{1}_m (\beta^{(i)})^T + (\zeta^{(i)}) \right] \boldsymbol{\theta},$$

where  $\mathbf{1}_m$  is the  $m$ -dimensional all-ones vector. By using matrix notation, we have managed to clearly express our desired payoff function as the linear combination of many function evaluations.

We now define the explore sampling distribution  $U : \Theta \times \mathcal{D} \rightarrow \Delta(\mathbb{R}^m \times \mathcal{C})$  as follows. With  $\frac{1}{4}$  probability, we return the point  $\mathbf{z}_{\text{exp}} = \mathbf{z}^{(i-1)}$  and weight vector  $\mathbf{w}_{\text{exp}} = (-2)\text{diag}(\mathbf{1}_m)\boldsymbol{\theta} = (-2)\mathbf{1}_m$ , where  $\text{diag}(\mathbf{1}_m)$  is the identity matrix with size  $m \times m$ . With  $\frac{1}{4}$  probability, we return the point  $\mathbf{z}_{\text{exp}} = \bar{\mathbf{z}}^{(i-1)}$  and weight vector  $\mathbf{w}_{\text{exp}} = (-2)\text{diag}(\mathbf{1}_m)\boldsymbol{\theta} = (-2)\mathbf{1}_m$ . For  $i = 1, \dots, m$ , with  $\frac{1}{4m}$  probability we return  $\mathbf{z}_{\text{exp}} = (\rho_i, \bar{\mathbf{z}}_{-i}^{(i-1)})$  and  $\mathbf{w}_{\text{exp}} = (4m) \left[ \frac{1}{2} \mathbf{1}_m \mathbf{e}_i^T + \text{diag}(\mathbf{e}_i) \mathbf{L}_{m,m} - \mathbf{L}_{m,m} \text{diag}(\mathbf{e}_i) \right] \boldsymbol{\theta}$ . For  $i = 1, \dots, m$ , with  $\frac{1}{4m}$  probability we return the point  $\mathbf{z}_{\text{exp}} = (\rho_i, \bar{\mathbf{z}}_{-i}^{(i-1)})$  and weight vector  $\mathbf{w}_{\text{exp}} = (4m) \left[ \frac{1}{2} \mathbf{1}_m \mathbf{e}_i^T + \text{diag}(\mathbf{e}_i) \mathbf{U}_{m,m} - \mathbf{U}_{m,m} \text{diag}(\mathbf{e}_i) \right] \boldsymbol{\theta}$ . Observe that, at subproblem  $i$  (essentially by construction):

$$\mathbb{E}[\hat{\mathbf{p}}(\boldsymbol{\theta}, \kappa)] = \mathbb{E}_{(\mathbf{z}^{(i-1)}, f) \sim \kappa} \mathbb{E}[\hat{\mathbf{p}}(\boldsymbol{\theta}, \text{AdvB}(\mathbf{z}^{(i-1)}, f))] \\ = \mathbb{E}_{(\mathbf{z}^{(i-1)}, f) \sim \kappa} \mathbb{E}_{(\mathbf{w}_{\text{exp}}, \mathbf{z}_{\text{exp}}) \sim \text{ExpS}(\boldsymbol{\theta}, \mathbf{z}^{(i-1)})} [f(\mathbf{z}_{\text{exp}}) \mathbf{w}_{\text{exp}}],$$

where

$$\begin{aligned}
& \mathbb{E}_{(\mathbf{w}_{\text{exp}}, \mathbf{z}_{\text{exp}}) \sim \text{ExpS}(\boldsymbol{\theta}, \mathbf{z}^{(i-1)})} [f(\mathbf{z}_{\text{exp}}) \mathbf{w}_{\text{exp}}] \\
&= \frac{1}{4} f(\mathbf{z}^{(i-1)}) [(-2) \text{diag}(\mathbf{1}_m) \boldsymbol{\theta}] + \frac{1}{4} f(\mathbf{z}^{(i-1)}) [(-2) \text{diag}(\mathbf{1}_m) \boldsymbol{\theta}] \\
&\quad + \sum_{i=1}^m \frac{1}{4m} f(\rho_i, \mathbf{z}_{-i}^{(i-1)}) \left[ (4m) \left[ \frac{1}{2} \mathbf{1}_m \mathbf{e}_i^T + \text{diag}(\mathbf{e}_i) \mathbf{L}_{m,m} - \mathbf{L}_{m,m} \text{diag}(\mathbf{e}_i) \right] \boldsymbol{\theta} \right] \\
&\quad + \sum_{i=1}^m \frac{1}{4m} f(\rho_i, \mathbf{z}_{-i}^{(i-1)}) \left[ (4m) \left[ \frac{1}{2} \mathbf{1}_m \mathbf{e}_i^T + \text{diag}(\mathbf{e}_i) \mathbf{U}_{m,m} - \mathbf{U}_{m,m} \text{diag}(\mathbf{e}_i) \right] \boldsymbol{\theta} \right] \\
&= f(\mathbf{z}^{(i-1)}) \left[ -\frac{1}{2} \mathbf{1}_m \mathbf{1}_m^T - \text{diag}(\mathbf{1}_m) \mathbf{L}_{m,m} + \mathbf{L}_{m,m} \text{diag}(\mathbf{1}_m) \right] \boldsymbol{\theta} \\
&\quad + f(\mathbf{z}^{(i-1)}) \left[ -\frac{1}{2} \mathbf{1}_m \mathbf{1}_m^T - \text{diag}(\mathbf{1}_m) \mathbf{U}_{m,m} + \mathbf{U}_{m,m} \text{diag}(\mathbf{1}_m) \right] \boldsymbol{\theta} \\
&\quad + \sum_{i=1}^m f(\rho_i, \mathbf{z}_{-i}^{(i-1)}) \left[ \left[ \frac{1}{2} \mathbf{1}_m \mathbf{e}_i^T + \text{diag}(\mathbf{e}_i) \mathbf{L}_{m,m} - \mathbf{L}_{m,m} \text{diag}(\mathbf{e}_i) \right] \boldsymbol{\theta} \right] \\
&\quad + \sum_{i=1}^m f(\rho_i, \mathbf{z}_{-i}^{(i-1)}) \left[ \left[ \frac{1}{2} \mathbf{1}_m \mathbf{e}_i^T + \text{diag}(\mathbf{e}_i) \mathbf{U}_{m,m} - \mathbf{U}_{m,m} \text{diag}(\mathbf{e}_i) \right] \boldsymbol{\theta} \right] \\
&= \left[ \frac{1}{2} \mathbf{1}_m (\boldsymbol{\alpha}^{(i)})^T + \frac{1}{2} \mathbf{1}_m (\boldsymbol{\beta}^{(i)})^T + \boldsymbol{\zeta}^{(i)} \right] \boldsymbol{\theta} \\
&= \text{PAYOFF}(\boldsymbol{\theta}, \mathbf{z}^{(i-1)}, f).
\end{aligned}$$

This explore sampling device also clearly runs in polynomial-time. Finally, we have

$$\begin{aligned}
\mathbb{E} [\hat{\mathbf{p}}(\boldsymbol{\theta}, \kappa)] &= \mathbb{E}_{(\mathbf{z}^{(i-1)}, f) \sim \kappa} \mathbb{E}_{(\mathbf{w}_{\text{exp}}, \mathbf{z}_{\text{exp}}) \sim \text{ExpS}(\boldsymbol{\theta}, \mathbf{z}^{(i-1)})} [f(\mathbf{z}_{\text{exp}}) \mathbf{w}_{\text{exp}}] \\
&= \mathbb{E}_{(\mathbf{z}^{(i-1)}, f) \sim \kappa} [\text{PAYOFF}(\boldsymbol{\theta}, \mathbf{z}^{(i-1)}, f)] = \mathbf{p}(\boldsymbol{\theta}, \kappa).
\end{aligned}$$

This completes the proof of bandit Blackwell reducibility.

For our bounds, we care about both the  $\ell_\infty$  diameter of the payoff  $\mathbf{D}(\mathbf{p})$  and the  $\ell_\infty$  diameter of the payoff estimator  $\mathbf{D}(\hat{\mathbf{p}})$ . The former is bounded by  $O(1)$ , since for any  $\boldsymbol{\theta}$ , the payoff function is a linear combination of  $O(1)$  function evaluations with  $O(1)$  coefficients. The latter is bounded by  $O(m)$  since aside from the  $O(4m)$ -scaling, the function evaluation yields a result in the range  $[0, 1]$  and the remaining terms have  $O(1)$  norms:

$$\begin{aligned}
\|\mathbf{1}_m\|_\infty &= 1 & \left\| \frac{1}{2} \mathbf{1}_m \mathbf{e}_i^T \boldsymbol{\theta} \right\|_\infty &= \frac{1}{2} [\boldsymbol{\theta}]_i \leq 1 \\
\|\text{diag}(\mathbf{e}_i) \mathbf{L}_{m,m} \boldsymbol{\theta}\|_\infty &= \sum_{j < i} [\boldsymbol{\theta}]_j \leq 1 & \|\mathbf{L}_{m,m} \text{diag}(\mathbf{e}_i)\|_\infty &= [\boldsymbol{\theta}]_i \leq 1 \\
\|\text{diag}(\mathbf{e}_i) \mathbf{U}_{m,m} \boldsymbol{\theta}\|_\infty &= \sum_{j > i} [\boldsymbol{\theta}]_j \leq 1 & \|\mathbf{U}_{m,m} \text{diag}(\mathbf{e}_i)\|_\infty &= [\boldsymbol{\theta}]_i \leq 1.
\end{aligned}$$

We complete the proof by applying Theorem 2 and Theorem 4, noting that our payoff dimension  $d$  equals the number of potential values that a coordinate can take,  $m$ :

$$\begin{aligned}
\frac{1}{2} \text{-regret}(\text{Algorithm 2 applied on Algorithm 7}) &\leq O(nT^{1/2} \log^{1/2} m) \\
\frac{1}{2} \text{-regret}(\text{Algorithm 4 applied on Algorithm 7}) &\leq O(nm^{2/3} T^{2/3} \log^{1/3} m).
\end{aligned}$$

■

## E.2. Proof of Corollaries 3 and 4

*Proof of Corollary 3.* We invoke Theorem 8 with the discretization space  $\mathcal{R} = \{0, 1\}$ . ■

*Proof of Corollary 4.* Let  $m \in \mathbb{Z}_+$  be a discretization parameter that we choose later to balance terms. We invoke Theorem 8 with discretization  $\mathcal{R} = \{0, \frac{1}{m}, \frac{2}{m}, \dots, 1\}$ . Because our functions are coordinate-wise  $L$ -Lipschitz, the (summed) discretization error is bounded by  $TL\frac{1}{m}n$  (note that the error from each subproblem is upper-bounded by  $L\frac{1}{m}$  and the error from each round is upper-bounded by  $L\frac{1}{m}n$ ). We choose  $m = LT^{1/2}$  for the full-information case and  $m = L^{3/5}T^{1/5}$  for the bandit case such that:

$$\begin{aligned} O(nT^{1/2} \log^{1/2} m) + TL\frac{1}{m}n &= O(nT^{1/2} \log^{1/2}(LT)) \\ O(nm^{2/3}T^{2/3} \log^{1/3} m) + TL\frac{1}{m}n &= O(nL^{2/5}T^{4/5} \log^{1/3}(LT)). \end{aligned}$$

This completes the proof. ■

## Appendix F: Application to Strong-DR Monotone Submodular Maximization over Downward Closed Convex Sets

*Preliminaries.* Consider the STRONG-DR MONOTONE SUBMODULAR MAXIMIZATION problem. Recall that a continuous multivariate function  $f : \mathbb{R}^n \rightarrow [0, 1]$  is strong-DR submodular if for all  $i \in [n]$ ,  $\mathbf{x} \preceq \mathbf{y} \in [0, 1]^n$ , and  $\delta \geq 0$ , we have

$$f(x_i + \delta, \mathbf{x}_{-i}) - f(\mathbf{x}) \geq f(y_i + \delta, \mathbf{y}_{-i}) - f(\mathbf{y}).$$

Here,  $\mathbf{x}_{-i}$  (resp.  $\mathbf{y}_{-i}$ ) is an  $(n-1)$ -dimensional vector with all coordinate values of  $\mathbf{x}$  (resp.  $\mathbf{y}$ ) except  $i$ , and  $\mathbf{x} \preceq \mathbf{y}$  if and only if  $\forall j \in [n] : x_j \leq y_j$ . Alternatively, a strong-DR continuous submodular function is a weak-DR continuous submodular function that is also concave along each coordinate. We also assume that function  $f : \mathbb{R}^n \rightarrow [0, 1]$  is  $L$ -Lipschitz smooth for some constant  $L > 0$ ; that is, for all  $\mathbf{x}, \mathbf{v} \in \mathbb{R}^n$ ,

$$f(\mathbf{x} + \mathbf{v}) - f(\mathbf{x}) \geq \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle - \frac{L}{2} \|\mathbf{v}\|_2^2,$$

where  $\nabla f(\mathbf{x}) = \left[ \frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_n} \right]^T$  for  $\mathbf{x} = [x_1, \dots, x_n]^T$ . Furthermore, we assume that the gradient of  $f$  has bounded  $\ell_\infty$  norm in some convex set  $\mathcal{P} \subseteq \mathbb{R}^n$ , that is, there is a constant  $U > 0$  such that  $\|\nabla f(\mathbf{x})\|_\infty \leq U$  for all  $\mathbf{x} \in \mathcal{P}$ . Note that as a simple corollary of this assumption, our functions are  $U$ -Lipschitz continuous in  $\ell_2$  norm, that is, for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ,

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq U \|\mathbf{x} - \mathbf{y}\|_2$$

### F.1. Offline Algorithm

Let  $\mathcal{P}$  be a downward closed convex polytope with  $\ell_2$  diameter  $R_{\mathcal{P}}$ , i.e.,  $R_{\mathcal{P}} = \max_{\mathbf{v}_1, \mathbf{v}_2 \in \mathcal{P}} \|\mathbf{v}_1 - \mathbf{v}_2\|_2$ . We further assume this convex set is polynomial-time separable; that is, we assume we have access to a polynomial-time algorithm for exactly solving any linear optimization over this convex set. More specifically, we assume that for any  $\mathbf{x} \in \mathcal{P}$ , we can solve the following optimization problem in polynomial time;  $\max_{\rho \in \mathcal{P}} \rho \cdot \nabla f(\mathbf{x})$ .

The goal of the offline problem is to maximize an  $L$ -Lipschitz smooth monotone strong-DR submodular functions  $f$  on  $\mathcal{P}$ , where  $\|\nabla f(\mathbf{z})\|_\infty \leq U$  for all  $\mathbf{z} \in \mathcal{P}$ . For the offline problem, the Frank-Wolfe variant for monotone strong-DR submodular function gives a  $(1 - 1/e)$ -approximation, which is known to be a tight approximation factor. This algorithm is first introduced in [Feldman et al. \(2011\)](#) and [Calinescu et al. \(2011\)](#) for the special case of multi-linear extension of set submodular functions, and is later extended to general strong-DR continuous submodular functions in [Bian et al. \(2017\)](#) and [Mokhtari et al. \(2020\)](#). We present this offline algorithm in [Algorithm 8](#). We then transform it into an online adversarial learning algorithm under the full-information ([Section F.2](#)) and bandit settings ([Section F.3](#)).

---

**Algorithm 8:** Frank-Wolfe variant for maximizing monotone strong-DR submodular function in a downward closed convex set. (c.f., [Bian et al. \(2017\)](#))

---

**Input:** Strong-DR monotone submodular function  $f : \mathcal{P} \rightarrow [0, 1]$  that is  $L$ -Lipschitz smooth with  $\|\nabla f(\mathbf{x})\|_\infty \leq U$  for all  $\mathbf{x} \in \mathcal{P}$ , where  $\mathcal{P}$  is a downward closed and convex polytope with  $\ell_2$  diameter  $R_{\mathcal{P}}$ . The number of iterations  $N$ .

**Output:** Point  $\mathbf{z} \in \mathcal{P}$ .

Initialize  $\mathbf{z}^{(0)} \leftarrow \mathbf{0}$ .

**for** iteration  $i = 1, 2, \dots, N$  **do**

Local Optimization Step

    Choose  $\boldsymbol{\rho}^{(i)} \in \mathcal{P}$  such that

$$\boldsymbol{\rho}^{(i)} \in \arg \max_{\boldsymbol{\rho} \in \mathcal{P}} \boldsymbol{\rho} \cdot \nabla f(\mathbf{z}^{(i-1)})$$

    (done in [Bian et al. \(2017\)](#) via maximizing a linear objective).

Local Update Step

    Set  $\mathbf{z}^{(i)} \leftarrow \mathbf{z}^{(i-1)} + \frac{1}{N} \boldsymbol{\rho}^{(i)}$ .

**return**  $\mathbf{z} \leftarrow \mathbf{z}^{(N)}$ .

---

The Frank-Wolfe variant algorithm takes the number of iterations  $N$  as input, which affects the approximation factor and also an additive error in its performance. The approximation factor of the algorithm for a finite  $N$  is  $\gamma_N = 1 - ((1 - 1/N))^N$ , which goes to  $1 - 1/e$  from above as  $N$  goes to infinity, and the additive error is in the order of  $O(\frac{LR_{\mathcal{P}}}{N})$ , which goes to zero as  $N$  goes to infinity (c.f., [Bian et al. \(2017\)](#)). That is,  $f(\mathbf{z}^{(N)}) \geq \gamma_N \max_{\mathbf{z} \in \mathcal{P}} f(\mathbf{z}) - \frac{LR_{\mathcal{P}}}{N}$ , where  $\mathbf{z}^{(N)}$  is the output of the Frank-Wolfe algorithm after  $N$  iterations. (We will also see the same bound later in the proof of [Theorem 9](#).) In the language of OFFLINE-IG, the parameter space is  $\Theta = \mathcal{P} \subset \mathbb{R}^n$  and  $d_{\text{param}} = n$ . In each subproblem  $i \in [N]$ , the algorithm picks the direction  $\boldsymbol{\rho}^{(i)}$  in the polytope  $\mathcal{P}$  that maximizes function  $\langle \boldsymbol{\rho}, \nabla f(\mathbf{z}^{(i-1)}) \rangle$ , (i.e.,  $\boldsymbol{\rho}^{(i)} \in \arg \max_{\boldsymbol{\rho} \in \mathcal{P}} \boldsymbol{\rho} \cdot \nabla f(\mathbf{z}^{(i-1)})$ .) Intuitively,  $\boldsymbol{\rho}^{(i)}$  is the direction along which we can maximize the improvement in the function value while still remaining feasible. Picking the direction inside the polytope eliminates the need for projecting the obtained

point back to  $\mathcal{P}$  at each iteration, which is usually an essential step in the Frank-Wolfe algorithm. We define the vector payoff function to be

$$\text{PAYOFF}(\boldsymbol{\rho}^{(i)}, \mathbf{z}^{(i-1)}, f) = \begin{bmatrix} -\boldsymbol{\rho}^{(i)} \cdot \mathbf{y}^{(i)} \\ \mathbf{y}^{(i)} \end{bmatrix}$$

where

$$\mathbf{y}^{(i)} \triangleq \nabla f(\mathbf{z}^{(i-1)}),$$

is the gradient of the function on the point from the previous iteration. The target set  $S$  is the polar cone of the  $\text{Cone}(1 \oplus \mathcal{P})$ , denoted by  $\text{Cone}(1 \oplus \mathcal{P})^\circ$ . Note that for  $\mathbf{x} \in \mathbb{R}^{d_1}$  and  $\mathbf{y} \in \mathbb{R}^{d_2}$ ,  $\mathbf{x} \oplus \mathbf{y} = \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \in \mathbb{R}^{d_1+d_2}$ . Moreover, for a cone  $C \subseteq \mathbb{R}^d$ , the polar cone of  $C$ , denoted by  $C^\circ$ , is defined as  $C^\circ = \{\boldsymbol{\theta} \in \mathbb{R}^d : \boldsymbol{\theta} \cdot \mathbf{x} \leq 0, \forall \mathbf{x} \in C\}$ . See that when  $\text{PAYOFF}(\boldsymbol{\rho}^{(i)}, \mathbf{z}^{(i-1)}, f) \in S$ , we have

$$-\boldsymbol{\rho}^{(i)} \cdot \mathbf{y}^{(i)} + \mathbf{y}^{(i)} \cdot \boldsymbol{\rho} \leq 0, \quad \forall \boldsymbol{\rho} \in \mathcal{P},$$

which implies

$$\boldsymbol{\rho}^{(i)} \in \arg \max_{\boldsymbol{\rho} \in \mathcal{P}} \boldsymbol{\rho} \cdot \mathbf{y}^{(i)} = \arg \max_{\boldsymbol{\rho} \in \mathcal{P}} \boldsymbol{\rho} \cdot \nabla f(\mathbf{z}^{(i-1)}).$$

So, picking  $\boldsymbol{\rho}^{(i)}$  that maximizes  $\boldsymbol{\rho} \cdot \nabla f(\mathbf{z}^{(i-1)})$  over  $\boldsymbol{\rho} \in \mathcal{P}$  is equivalent to picking  $\boldsymbol{\rho}^{(i)}$  such that the payoff  $\text{PAYOFF}(\boldsymbol{\rho}^{(i)}, \mathbf{z}^{(i-1)}, f)$  is in  $S$ .

Algorithm 8 is a variant of OFFLINE-IG (Algorithm 1) where the local optimization step is replaced with picking  $\boldsymbol{\rho}^{(i)} \in \mathcal{P}$  such that the payoff function evaluated on  $\boldsymbol{\rho}^{(i)}$  falls in the target set  $S$ . Recall that in OFFLINE-IG,  $\boldsymbol{\rho}^{(i)}$  is picked in the local optimization step such that the payoff function on  $\boldsymbol{\rho}^{(i)}$  is non-negative, i.e.,  $\text{PAYOFF}(\boldsymbol{\rho}^{(i)}, \mathbf{z}^{(i-1)}, f) \geq \mathbf{0}$ . In this variant, we want to pick  $\boldsymbol{\rho}^{(i)}$  such that  $\text{PAYOFF}(\boldsymbol{\rho}^{(i)}, \mathbf{z}^{(i-1)}, f) \in S$  for some convex set  $S$ , which is not a positive orthant.

Now, consider a Blackwell sequential game, as defined in Section 2.3, where  $\mathcal{X} = \Theta = \mathcal{P}$ ,  $\mathcal{Y} = [-U, U]^n$ ,  $\mathbf{p}(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} -\mathbf{x} \cdot \mathbf{y} \\ \mathbf{y} \end{bmatrix} \in \mathbb{R}^{n+1}$ , and  $S = \text{Cone}(1 \oplus \mathcal{P})^\circ \subset \mathbb{R}^{n+1}$ . Note that the target set and payoff are  $(n+1)$ -dimensional. The set  $S$  is response-satisfiable, since for every player 2's action  $\mathbf{y} \in \mathcal{Y}$ , there exists a player 1's action,  $\mathbf{x}^* \in \arg \max_{\mathbf{x} \in \mathcal{P}} \mathbf{x} \cdot \mathbf{y}$ , such that  $\mathbf{p}(\mathbf{x}^*, \mathbf{y}) \in S$ . Specifically, when  $\mathbf{x}^* \in \arg \max_{\mathbf{x} \in \mathcal{P}} \mathbf{x} \cdot \mathbf{y}$ , then  $\mathbf{x}^* \cdot \mathbf{y} \geq \mathbf{x} \cdot \mathbf{y}$  for all  $\mathbf{x} \in \mathcal{P}$ , which implies  $\mathbf{p}(\mathbf{x}^*, \mathbf{y}) \cdot [1, \mathbf{x}]^T \leq 0$  for all  $\mathbf{x} \in \mathcal{P}$ , further implying  $\mathbf{p}(\mathbf{x}^*, \mathbf{y}) \in S$ . Note that since  $S$  is response-satisfiable and  $S$  is polynomial-time separable, there exists a polynomial-time algorithm AlgB, based on Theorem 1, such that for any sequence of actions  $\mathbf{y}_1, \dots, \mathbf{y}_T \in \mathcal{Y}$  generated by an adaptive adversary, AlgB can generate a sequence of actions  $\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathcal{P}$  such that

$$d_\infty \left( \frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S \right) \leq O \left( D_\infty(\mathcal{P}) \sqrt{\frac{\log d_{\mathcal{P}}}{T}} \right) = O \left( U \sqrt{\frac{\log n}{T}} \right). \quad (23)$$

## F.2. Online Learning Algorithm for Full-Information Setting

In the online problem, we consider an adversarial sequence of strong-DR submodular functions  $f_1, \dots, f_T$ , where for each  $1 \leq t \leq T$ ,  $f_t: \mathcal{P} \rightarrow [0, 1]$  for some downward closed polynomial-time separable convex polytope  $\mathcal{P} \subseteq \mathbb{R}^n$  with  $\ell_2$  diameter  $R_{\mathcal{P}}$ . Moreover, for each  $t \in [T]$ ,  $f_t$  is  $L$ -Lipschitz smooth and  $\|\nabla f_t(\mathbf{z})\|_{\infty} \leq U$  for all  $\mathbf{z} \in \mathcal{P}$ . Then, for each subproblem  $i \in [N]$  in Algorithm 8, we can replace the local optimization step with a polynomial-time online algorithm  $\text{AlgB}^{(i)}$  where  $\mathbf{x}_t = \boldsymbol{\rho}_t^{(i)}$  and  $\mathbf{y}_t = \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t) = \nabla f_t(\mathbf{z}_t^{(i-1)})$ . For subproblem  $i$ , the guarantee in Equation (23) becomes

$$d_{\infty} \left( \frac{1}{T} \sum_{t=1}^T \text{PAYOFF}(\boldsymbol{\rho}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t), S \right) \leq O \left( U \sqrt{\frac{\log n}{T}} \right). \quad (24)$$

We utilize  $N$  parallel runs of  $\text{AlgB}$ , one for each subproblem as shown in Algorithm 9. Note that we do not directly apply our meta Algorithm 2 in this problem because the target set  $S$  is not the positive orthant. But due to this subtle difference, we derive our results with similar approach and our algorithm follows the same framework.

---

### Algorithm 9: Full-information online learning algorithm for Algorithm 8

---

**Input:** A sequence of  $L$ -Lipschitz smooth strong-DR monotone submodular functions

$f_1, \dots, f_T: \mathcal{P} \rightarrow [0, 1]$  such that  $\|\nabla f_t(\mathbf{z})\|_{\infty} \leq U$  for all  $\mathbf{z} \in \mathcal{P}, t \in [T]$ , where  $\mathcal{P}$  is a downward closed and convex polytope with  $\ell_2$  diameter  $R_{\mathcal{P}}$ .  $N$  algorithms  $\{\text{AlgB}^{(i)}\}_{i=1, \dots, N}$ , where each  $\text{AlgB}^{(i)}$  is an online learning algorithm for a Blackwell sequential game with  $\mathcal{X} = \mathcal{P}, \mathcal{Y} = [-U, U]^n, \mathbf{p}(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} -\mathbf{x} \cdot \mathbf{y} \\ \mathbf{y} \end{bmatrix}$ , and  $S = \text{Cone}(1 \oplus \mathcal{P})^{\circ}$ .

**Output:**  $\mathbf{z}_1, \dots, \mathbf{z}_T \in \mathcal{P}$ .

**for**  $t = 1, \dots, T$  **do**

Initialize  $\mathbf{z}_t^{(0)} \leftarrow \mathbf{0}$ .

**for** iteration  $i = 1, 2, \dots, N$  **do**

Choose  $\boldsymbol{\rho}_t^{(i)} \in \mathcal{P}$  by querying online algorithm  $\text{AlgB}^{(i)}$  given the update parameters and vector payoffs prior to round  $t$  in the Blackwell sequential game of subproblem  $i$ , i.e.,  $\boldsymbol{\rho}_t^{(i)} \leftarrow \text{AlgB}^{(i)} \left( \boldsymbol{\rho}_1^{(i)}, \dots, \boldsymbol{\rho}_{t-1}^{(i)}, \{\mathbf{p}(\boldsymbol{\rho}_{\tau}^{(i)}, \mathbf{y}_{\tau}^{(i)})\}_{\tau \in [1:t-1]} \right)$ .

Set  $\mathbf{z}_t^{(i)} \leftarrow \mathbf{z}_t^{(i-1)} + \frac{1}{N} \boldsymbol{\rho}_t^{(i)}$ .

**return**  $\mathbf{z}_t \leftarrow \mathbf{z}_t^{(N)}$ .

*Adversary reveals function  $f_t$ .*

**for** iteration  $i = 1, 2, \dots, N$  **do**

Compute  $\nabla f_t(\mathbf{z}_t^{(i-1)})$ .

Give feedback  $\mathbf{p}(\boldsymbol{\rho}_t^{(i)}, \mathbf{y}_t^{(i)}) \leftarrow \begin{bmatrix} -\boldsymbol{\rho}_t^{(i)} \cdot \nabla f_t(\mathbf{z}_t^{(i-1)}) \\ \nabla f_t(\mathbf{z}_t^{(i-1)}) \end{bmatrix}$  to the Blackwell algorithm  $\text{AlgB}^{(i)}$ .

Note that  $\mathbf{y}_t^{(i)} = \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t) = \nabla f_t(\mathbf{z}_t^{(i-1)})$ .

---

**THEOREM 9. (Online learning for strong-DR monotone SM maximization over downward closed convex sets)** Let  $\mathcal{P}$  be a downward closed convex polytope with  $\ell_2$  diameter  $R_{\mathcal{P}}$ . Consider the online learning problem of maximizing strong-DR monotone submodular functions on  $\mathcal{P}$ , which are  $L$ -Lipschitz smooth and have bounded gradient of  $U$  for any point  $\mathbf{z} \in \mathcal{P}$ . Then, under full-information setting, Algorithm 9 with  $N$  iterations obtains  $O\left(R_{\mathcal{P}}U\sqrt{Tn\log n} + \frac{TLR_{\mathcal{P}}^2}{N}\right) \gamma_N$ -regret, where  $T$  is the number of rounds and  $\gamma_N = 1 - (1 - 1/N)^N \geq (1 - 1/e)$ . The benchmark in the regret bounds is  $\gamma_N \max_{\mathbf{z} \in \mathcal{P}} \sum_{t=1}^T f_t(\mathbf{z})$ . By setting  $N = \frac{\sqrt{T}LR_{\mathcal{P}}}{U\sqrt{n\log n}}$ , we get  $O\left(R_{\mathcal{P}}U\sqrt{Tn\log n}\right) \gamma_N$ -regret.

Before proving the above theorem, we first prove the following technical lemma using strong-DR submodularity property of our functions. The proof of Lemma 5 is presented at the end of this section.

**LEMMA 5.** Consider any iteration  $i \in [N]$  of the Frank-Wolfe algorithm (Algorithm 8). For all  $\mathbf{z} \in \mathcal{P}$ , we have

$$\boldsymbol{\rho}^{(i)} \cdot \nabla f(\mathbf{z}^{(i-1)}) \geq f(\mathbf{z}) - f(\mathbf{z}^{(i-1)}).$$

*Proof of Theorem 9.* We show that Algorithm 9 works in two steps: (i) proving a variant of the extended robustness property and (ii) proving a variant of the implication of Blackwell reducibility.

We start with a variant of the extended robustness property. Suppose that  $\{\boldsymbol{\rho}_t^{(i)}, \mathbf{z}_t^{(i)}\}_{t \in [T], 0 \leq i \leq N}$  are the values from Algorithm 9 with  $N$  iterations. We show that if the following equation holds for some function  $h$ :

$$\forall i \in [N], \quad \sum_{t=1}^T \boldsymbol{\rho}_t^{(i)} \cdot \nabla f_t(\mathbf{z}_t^{(i-1)}) \geq \max_{\boldsymbol{\rho} \in \mathcal{P}} \sum_{t=1}^T \boldsymbol{\rho} \cdot \nabla f_t(\mathbf{z}_t^{(i-1)}) - h(T), \quad (25)$$

then we have

$$\forall \mathbf{z}^* \in \mathcal{P}, \quad \sum_{t=1}^T f_t(\mathbf{z}_t) \geq \gamma_N \sum_{t=1}^T f_t(\mathbf{z}^*) - h(T) - \frac{TLR_{\mathcal{P}}^2}{2N}, \quad (26)$$

where  $\gamma_N = 1 - (1 - \frac{1}{N})^N$ . To show this, we use Lemma 5. Note that for each round  $t$  and subproblem  $i \in [N]$ , we have

$$\begin{aligned} f_t(\mathbf{z}_t^{(i)}) - f_t(\mathbf{z}_t^{(i-1)}) &= f_t(\mathbf{z}_t^{(i-1)}) + \frac{1}{N} \boldsymbol{\rho}_t^{(i)} \cdot \nabla f_t(\mathbf{z}_t^{(i-1)}) - f_t(\mathbf{z}_t^{(i-1)}) \\ &\stackrel{(1)}{\geq} \frac{1}{N} \nabla f_t(\mathbf{z}_t^{(i-1)}) \cdot \boldsymbol{\rho}_t^{(i)} - \frac{1}{2} \frac{L \|\boldsymbol{\rho}_t^{(i)}\|_2^2}{N^2} \\ &\stackrel{(2)}{\geq} \frac{1}{N} \nabla f_t(\mathbf{z}_t^{(i-1)}) \cdot \boldsymbol{\rho}_t^{(i)} - \frac{1}{2} \frac{LR_{\mathcal{P}}^2}{N^2}. \end{aligned}$$

Inequality (1) holds because  $f_t$  is  $L$ -Lipschitz smooth and Inequality (2) holds since  $\boldsymbol{\rho}_t^{(i)} \in \mathcal{P}$  and  $\mathcal{P}$  is a downward closed convex polytope with  $\ell_2$  diameter  $R_{\mathcal{P}}$ . This implies that  $\|\boldsymbol{\rho}_t^{(i)}\|_2^2 \leq R_{\mathcal{P}}^2$ .

Summing this for all  $t \in [T]$ , for any  $\mathbf{z}^* \in \mathcal{P}$ , we have

$$\begin{aligned} \sum_{t=1}^T \left( f_t(\mathbf{z}_t^{(i)}) - f_t(\mathbf{z}_t^{(i-1)}) \right) &\geq \frac{1}{N} \sum_{t=1}^T \nabla f_t(\mathbf{z}_t^{(i-1)}) \cdot \boldsymbol{\rho}_t^{(i)} - \frac{1}{2} \frac{TLR_{\mathcal{P}}^2}{N^2} \\ &\stackrel{(1)}{\geq} \frac{1}{N} \max_{\boldsymbol{\rho} \in \mathcal{P}} \sum_{t=1}^T \boldsymbol{\rho} \cdot \nabla f_t(\mathbf{z}_t^{(i-1)}) - \frac{1}{N} h(T) - \frac{1}{2} \frac{TLR_{\mathcal{P}}^2}{N^2} \\ &\stackrel{(2)}{\geq} \frac{1}{N} \left( \sum_{t=1}^T f_t(\mathbf{z}^*) - \sum_{t=1}^T f_t(\mathbf{z}_t^{(i-1)}) \right) - \frac{1}{N} h(T) - \frac{1}{2} \frac{TLR_{\mathcal{P}}^2}{N^2}, \end{aligned}$$

where Inequality (1) holds because of the assumption in Equation (25), and Inequality (2) is reached by applying Lemma 5. Rearranging the terms, for any  $\mathbf{z}^* \in \mathcal{P}$ , we have

$$\sum_{t=1}^T \left( f_t(\mathbf{z}_t^{(i)}) - f_t(\mathbf{z}^*) \right) \geq \left( 1 - \frac{1}{N} \right) \left( \sum_{t=1}^T \left( f_t(\mathbf{z}_t^{(i-1)}) - f_t(\mathbf{z}^*) \right) \right) - \frac{1}{N} h(T) - \frac{1}{2} \frac{TLR_{\mathcal{P}}^2}{N^2},$$

and iterating from  $i = 1, 2, \dots, N$ , we get

$$\sum_{t=1}^T \left( f_t(\mathbf{z}_t^{(N)}) - f_t(\mathbf{z}^*) \right) \geq \left( 1 - \frac{1}{N} \right)^N \left( \sum_{t=1}^T \left( f_t(\mathbf{z}_t^{(0)}) - f_t(\mathbf{z}^*) \right) \right) - h(T) - \frac{1}{2} \frac{TLR_{\mathcal{P}}^2}{N}.$$

Since  $f_t(\mathbf{z}_t^{(0)}) = f_t(\mathbf{0}) \geq 0$ , for any  $\mathbf{z}^* \in \mathcal{P}$ , we have

$$\sum_{t=1}^T f_t(\mathbf{z}_t^{(N)}) \geq \left( 1 - \left( 1 - \frac{1}{N} \right)^N \right) \sum_{t=1}^T f_t(\mathbf{z}^*) - h(T) - \frac{TLR_{\mathcal{P}}^2}{2N} = \gamma_N \sum_{t=1}^T f_t(\mathbf{z}^*) - h(T) - \frac{TLR_{\mathcal{P}}^2}{2N},$$

which is the desired inequality in Equation (26).

We now prove that the regret guarantee in Equation (24), which is obtained from Theorem 1, implies Equation (25) for  $h(T) = O(R_{\mathcal{P}}U\sqrt{Tn\log n})$ . From the Blackwell guarantee (i.e., Equation (24)), we have

$$O\left(U\sqrt{\frac{\log n}{T}}\right) \geq d_{\infty} \left( \frac{1}{T} \sum_{t=1}^T \text{PAYOFF}(\boldsymbol{\rho}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t), S \right) \quad (27)$$

$$\stackrel{(1)}{\geq} \frac{1}{\sqrt{n+1}} d_2 \left( \frac{1}{T} \sum_{t=1}^T \text{PAYOFF}(\boldsymbol{\rho}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t), S \right)$$

$$\stackrel{(2)}{=} \frac{1}{\sqrt{n+1}} \max_{\mathbf{w} \in \text{Cone}(1 \oplus \mathcal{P}) \cap B_2^{n+1}(1)} \mathbf{w} \cdot \frac{1}{T} \sum_{t=1}^T \text{PAYOFF}(\boldsymbol{\rho}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t)$$

$$\stackrel{(3)}{=} \frac{1}{\sqrt{n+1}} \max_{\boldsymbol{\rho} \in \mathcal{P}} \frac{(1 \oplus \boldsymbol{\rho})}{\|1 \oplus \boldsymbol{\rho}\|_2} \cdot \frac{1}{T} \sum_{t=1}^T \text{PAYOFF}(\boldsymbol{\rho}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t)$$

$$\stackrel{(4)}{=} \frac{1}{\sqrt{n+1}} \max_{\boldsymbol{\rho} \in \mathcal{P}} \frac{\left( \frac{1}{T} \sum_{t=1}^T \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho} - \frac{1}{T} \sum_{t=1}^T \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}_t^{(i)} \right)}{\|1 \oplus \boldsymbol{\rho}\|_2}$$

$$\stackrel{(5)}{\geq} \frac{1}{\sqrt{n+1}\sqrt{1+R_{\mathcal{P}}^2}} \frac{1}{T} \left( \max_{\boldsymbol{\rho} \in \mathcal{P}} \sum_{t=1}^T \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho} - \sum_{t=1}^T \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}_t^{(i)} \right) \quad (28)$$

Here,  $B_2^{n+1}(1)$  is a  $(n+1)$ -dimensional  $\ell_2$  unit ball. Inequality (1) holds since the dimension of the payoff vector is  $(n+1)$ , Inequality (2) follows from Lemma 13 in Abernethy et al. (2011). (The lemma says that for every convex cone  $C$  in  $\mathbb{R}^d$ ,  $d_2(\mathbf{x}, C) = \max_{\boldsymbol{\theta} \in C \cap B_2^d(1)} \langle \boldsymbol{\theta}, \mathbf{x} \rangle$ .) Inequality (3) follows from rewriting  $\mathbf{w}$  as a function of  $\boldsymbol{\rho} \in \mathcal{P}$ , Inequality (4) holds because the definition of the payoff vector where  $\mathbf{y}_t^{(i)} = \nabla f_t(\mathbf{z}_t^{(i-1)})$ , and Inequality (5) holds because  $\|1 \oplus \boldsymbol{\rho}\|_2 \leq \sqrt{1+R_{\mathcal{P}}^2}$  for any  $\boldsymbol{\rho} \in \mathcal{P}$ . Putting things together and rearranging, we get

$$\sum_{t=1}^T \boldsymbol{\rho}_t^{(i)} \cdot \mathbf{y}_t^{(i)} \geq \max_{\boldsymbol{\rho} \in \mathcal{P}} \boldsymbol{\rho} \cdot \sum_{t=1}^T \mathbf{y}_t^{(i)} - O\left(U\sqrt{\frac{\log n}{T}} \cdot R_{\mathcal{P}}\sqrt{nT}\right) = \max_{\boldsymbol{\rho} \in \mathcal{P}} \boldsymbol{\rho} \cdot \sum_{t=1}^T \mathbf{y}_t^{(i)} - O\left(R_{\mathcal{P}}U\sqrt{Tn\log n}\right),$$

which is Equation (25) with  $h(T) = O(R_{\mathcal{P}}U\sqrt{Tn\log n})$ .

In conclusion, by invoking Equation (26) with  $h(T) = O(R_{\mathcal{P}}U\sqrt{Tn\log n})$ , the total regret of Algorithm 9 is  $O\left(R_{\mathcal{P}}U\sqrt{Tn\log n} + \frac{TLR_{\mathcal{P}}^2}{N}\right)$ . By setting  $N = \frac{\sqrt{T}LR_{\mathcal{P}}}{U\sqrt{n\log n}}$ , we get  $O\left(R_{\mathcal{P}}U\sqrt{Tn\log n}\right)$   $\gamma_N$ -regret, as desired. ■

For completeness, we provide the proof of the technical lemma used in the above proof.

*Proof of Lemma 5.* Let  $\boldsymbol{\rho}^* = \mathbf{z} \vee \mathbf{z}^{(i-1)} - \mathbf{z}^{(i-1)}$ . Observe that

$$\begin{aligned} f(\mathbf{z}) - f(\mathbf{z}^{(i-1)}) &\stackrel{(1)}{\leq} f(\mathbf{z} \vee \mathbf{z}^{(i-1)}) - f(\mathbf{z}^{(i-1)}) = f(\mathbf{z}^{(i-1)} + \boldsymbol{\rho}^*) - f(\mathbf{z}^{(i-1)}) \\ &\stackrel{(2)}{\leq} \boldsymbol{\rho}^* \cdot \nabla f(\mathbf{z}^{(i-1)}) \stackrel{(3)}{\leq} \boldsymbol{\rho}^{(i)} \cdot \nabla f(\mathbf{z}^{(i-1)}), \end{aligned}$$

where Inequality (1) holds because  $f$  is a monotone function, Inequality (2) holds because of the concavity (along each direction) property of strong-DR submodular functions, and Inequality (3) holds because  $\boldsymbol{\rho}^{(i)}$  maximizes  $\boldsymbol{\rho} \cdot \nabla f(\mathbf{z}^{(i-1)})$  over  $\mathcal{P}$  and  $\boldsymbol{\rho}^* \in \mathcal{P}$  since  $\mathcal{P}$  is downward closed. ■

**F.2.1. Application to Set Submodular Maximization Subject to Matroids** As an important application of Theorem 9, we consider the online learning version of maximizing monotone set submodular functions subject to a matroid constraint. We start by stating the following simple corollary of Theorem 9 for maximizing strong-DR monotone submodular maximization over matroid polytopes, which is the key to obtain similar results for maximizing monotone set submodular maximization subject to matroids constraints.

**COROLLARY 5.** *Let  $M = (V, \mathcal{I})$  be a matroid with  $|V| = n$  and  $\text{Rank}(M) = R$ . Define  $\mathcal{P}_M = \text{Conv}(\{\mathbf{1}_I : I \in \mathcal{I}\})$  as a matroid polytope associated with matroid  $M = (V, \mathcal{I})$ . When we restrict our domain to  $\mathcal{P}_M$  (i.e., if we set  $\mathcal{P}$  to  $\mathcal{P}_M$ ), we obtain a total regret of  $O\left(RU\sqrt{Tn\log n} + \frac{TLR^2}{Nn}\right)$  for Algorithm 9. Setting  $N = \frac{\sqrt{TLR}}{Un\sqrt{n\log n}}$ , we get  $O\left(RU\sqrt{Tn\log n}\right) \gamma_N$ -regret.*

*Proof of Corollary 5.* First observe that the  $\ell_2$  diameter of the polytope  $\mathcal{P}_M$  becomes  $\frac{R}{\sqrt{n}}$  since for any  $\boldsymbol{\rho}$  in matroid polytope  $\mathcal{P}_M$ , we have

$$\|\boldsymbol{\rho}\|_2^2 \leq n \left( \frac{\text{Rank}(M)}{n} \right)^2 = \frac{R^2}{n}.$$

Moreover, we have  $\|\mathbf{1} \oplus \boldsymbol{\rho}\|_2 \leq \max\{1, \sqrt{2} \max_{\boldsymbol{\rho} \in \mathcal{P}} \|\boldsymbol{\rho}\|_2\} \leq \max\{1, \sqrt{2} \frac{R}{\sqrt{n}}\}$ , which results in  $h(T) = O\left(RU\sqrt{Tn\log n}\right)$  in Equation (25). This implies a total regret of  $O\left(RU\sqrt{Tn\log n} + \frac{TLR^2}{Nn}\right)$  for Algorithm 9. Setting  $N = \frac{\sqrt{TLR}}{Un\sqrt{n\log n}}$ , we get  $O\left(RU\sqrt{Tn\log n}\right) \gamma_N$ -regret.

Next, we seek to understand if we can obtain similar results as in Corollary 5 for maximizing monotone set submodular functions with a matroid constraint. To do so, we rely on the notion of *multi-linear extension* of set submodular functions and *pipage rounding algorithm* proposed in Calinescu et al. (2011). Using these two notions, Calinescu et al. (2011) further propose an optimal approximation algorithm for maximizing monotone set submodular functions with a matroid constraint. We briefly describe this approach below and then show how to combine it with Algorithm 8 to extend our result to the full information online learning version of this problem.<sup>28</sup>

Given a set submodular function  $f : 2^{[n]} \rightarrow [0, 1]$ , the following continuous extension of this function, known as the multi-linear extension, is a strong-DR submodular continuous function (Calinescu et al. 2011):

$$\forall \mathbf{x} \in [0, 1]^n : f^{\text{MLE}}(\mathbf{x}) = \sum_{S \subseteq [n]} f(S) \prod_{i \in S} x_i \prod_{i \notin S} (1 - x_i). \quad (29)$$

<sup>28</sup> Because of a subtle technical reason, this result does not immediately extend to the bandit setting; see Remark 6.

Notably,  $f^{\text{MLE}}(\mathbf{x})$  is essentially the expected value of  $f(S)$  at a randomized set  $S$ , where each element  $i \in [n]$  is placed in  $S$  independently with probability  $x_i$ . As a result, (i) the two functions have the same value when  $\mathbf{x} \in \{0, 1\}^n$ , and (ii) if  $\mathbf{x}^*$  maximizes  $f(\mathbf{x})$  over a matroid polytope,  $f^{\text{MLE}}(\mathbf{x}^*)$  should be equal to the maximum value of  $f(S)$  over independent sets  $S$  of its associated matroid. This last property is a simple consequence of the existence of polynomial-time loss-less randomized rounding algorithms in the matroid polytope for monotone submodular functions, such as the pipage rounding introduced in Calinescu et al. (2011). Given a point  $\mathbf{x}$  in the matroid polytope, the pipage randomized rounding algorithm returns a randomized set  $\tilde{S}$  such that  $\tilde{S}$  is always an independent set of the matroid and  $\mathbb{E}[f(\tilde{S})] \geq f^{\text{MLE}}(\mathbf{x})$ . Applying such a loss-less rounding algorithm to the maximizer point  $\mathbf{x}^*$  will result in a distribution over independent sets of the matroid, where every point in the support of this distribution should be a maximizer of  $f(S)$  over independent sets of the polytope (otherwise,  $\mathbf{x}^*$  could not be the maximizer). We should also highlight that one can aim to approximately maximize the multi-linear extension of the monotone submodular function subject to a matroid constraint as it is strong-DR, Lipschitz continuous and smooth, e.g., see Bian et al. (2017). This can be done using a first-order method such as Algorithm 8. Then, we can use the pipage rounding algorithm to obtain a randomized approximation algorithm with *exactly* the same approximation guarantee in-expectation for monotone set submodular maximization subject to a matroid constraint.

We now sketch how to use the multi-linear extension and the pipage rounding algorithm in a full-information online learning setting, and how to obtain a similar result as in Corollary 5 but for monotone set submodular functions  $f_1, \dots, f_T$ . The idea is running the online learning algorithm of the previous section (Algorithm 9) on the sequence of functions  $f_1^{\text{MLE}}, \dots, f_T^{\text{MLE}}$  (which are strong-DR continuous submodular functions) given the full information feedback  $f_t(\cdot)$  at the end of each round  $t$  (will be detailed later). Running Algorithm 9 generates points  $\mathbf{z}_1, \dots, \mathbf{z}_T$  in the matroid polytope. We then sample randomized sets  $\tilde{S}_1, \dots, \tilde{S}_T$  at the end of each round  $t$  using the pipage randomized rounding algorithms, so that these sets are independent sets in the matroid and also  $\mathbb{E}[f_t(\tilde{S}_t)] \geq f^{\text{MLE}}(\mathbf{z}_t)$ . In this way, we have:

$$\gamma_N \cdot \max_{S \in \mathcal{I}} \sum_{t \in [T]} f_t(S) - \sum_{t \in [T]} \mathbb{E}[f_t(\tilde{S}_t)] \leq \gamma_N \cdot \max_{\mathbf{z} \in \text{Conv}(\mathcal{I})} \sum_{t \in [T]} f_t^{\text{MLE}}(\mathbf{z}) - \sum_{t \in [T]} f_t^{\text{MLE}}(\mathbf{z}_t),$$

and therefore the regret bound of Algorithm 9 carries over to this setting.

To run the full-information online learning algorithm Algorithm 9 as described above, we need to be able to construct the full information feedback  $f_t^{\text{MLS}}(\mathbf{z})$  (for all  $\mathbf{z} \in [0, 1]^n$ ) for the underlying strong-DR continuous submodular function  $f_t^{\text{MLE}}$  given the set function full-information feedback  $f_t$  that we receive in each round. Given the full information feedback  $f_t$ ,  $f_t^{\text{MLS}}(\mathbf{z})$  can be estimated at any point  $\mathbf{z}$  through sampling and taking average. The sampling procedure evaluates  $f_t$  at a randomized set  $S$ , where every element  $i$  is placed in  $S$  independently with probability  $z_i$ . By repeating the same process and averaging one can obtain an estimation of  $f_t^{\text{MLS}}(\mathbf{z})$  at any accuracy level. (For example, with samples complexity of  $\Omega(T^4)$ , the accuracy level in estimating multi-linear extension (and hence its gradient) at any point will be of  $O(\frac{1}{T^2})$  by applying a concentration bound such as the Chernoff-Hoeffding bound. Note that such an additive error is negligible in our regret bounds as it causes an additional regret that is in the order of  $o(1)$  by following similar lines as in the proof of Theorem 9. We omit the details for brevity.)

### F.3. Online Learning Algorithm for Bandit Setting

In the bandit setting, we observe the function value evaluated on the chosen point at each round, instead of getting the full function  $f_t$ . That prevents us from providing the right feedback for each of the Blackwell algorithms. Recall that in Algorithm 8, the algorithm gives the feedback

$$\mathbf{p}(\boldsymbol{\rho}_t^{(i)}, \mathbf{y}_t^{(i)}) \leftarrow \begin{bmatrix} -\boldsymbol{\rho}_t^{(i)} \cdot \nabla f_t(\mathbf{z}_t^{(i-1)}) \\ \nabla f_t(\mathbf{z}_t^{(i-1)}) \end{bmatrix}$$

to  $\text{AlgB}^{(i)}$  in each round. Given this, to design an online learning algorithm for the bandit setting, we would like to provide an unbiased estimator of the gradient  $\nabla f_t(\mathbf{z}_t^{(i-1)})$  in terms of the function value  $f_t(\mathbf{z}_t^{(i-1)})$  every time the algorithm explores. To do so, following Zhang et al. (2019), we use an intermediate function  $\widehat{f}_{\delta,t}(\mathbf{z}) = \mathbb{E}_{\mathbf{v} \sim \text{Unif}(B_2^n)}[f(\mathbf{z} + \delta\mathbf{v})]$  for some small value of  $\delta > 0$  that we will determine later, where  $B_2^n$  denotes the  $\ell_2$  unit ball in  $\mathbb{R}^n$ . Note that when  $f_t$  is monotone, continuous strong-DR submodular,  $L$ -Lipschitz smooth with a bounded gradient of  $U$  for any point  $\mathbf{z} \in \mathcal{P}$ , so is  $\widehat{f}_{\delta,t}$ . We further have  $|f_t(\mathbf{z}) - \widehat{f}_{\delta,t}(\mathbf{z})| \leq U\delta$  for all  $\mathbf{z} \in \mathcal{P}$ . Moreover, let  $S_2^{(n-1)}$  denote the  $\ell_2$  unit sphere in  $\mathbb{R}^n$  (i.e.,  $S_2^{(n-1)} = \{\mathbf{v} \in \mathbb{R}^n : \|\mathbf{v}\|_2 = 1\}$ ). Then, when  $\mathbf{u}$  is uniformly drawn at random from  $S_2^{(n-1)}$ ,  $\frac{n}{\delta} f_t(\mathbf{z} + \delta\mathbf{u})\mathbf{u}$  is an unbiased estimator of  $\nabla \widehat{f}_{\delta,t}(\mathbf{z})$ , i.e.,  $\nabla \widehat{f}_{\delta,t}(\mathbf{z}) = \mathbb{E}_{\mathbf{u} \sim \text{Unif}(S_2^{(n-1)})} \left[ \frac{n}{\delta} f_t(\mathbf{z} + \delta\mathbf{u})\mathbf{u} \right]$  (Flaxman et al. (2005), Zhang et al. (2019)). We further use an intermediate convex downward closed feasible region  $\mathcal{P}' = (1 - \alpha)\mathcal{P} + \delta\mathbf{1}$ , where  $\alpha = \frac{\sqrt{n+1}}{r}\delta$ , and  $r$  is the largest positive real number such that  $r \cdot (B_2^n \cap \mathbb{R}_{\geq 0}^n) \subseteq \mathcal{P}$ . The new feasible region  $\mathcal{P}'$  is such that for any  $\mathbf{z} \in \mathcal{P}'$  and  $\mathbf{u} \in S^{(n-1)}$ , the point  $\mathbf{z} + \delta\mathbf{u}$  is in  $\mathcal{P}$ , keeping the chosen point of  $\mathbf{z} + \delta\mathbf{u}$  feasible in the original problem.

In the bandit algorithm, after showing bandit Blackwell reducibility as in Theorem 3, we replace the local optimization step with a polynomial-time online algorithm  $\text{AlgBB}^{(i)}$  for each subproblem  $i \in [N]$ . We utilize  $N$  parallel runs of  $\text{AlgBB}$ , one for each subproblem as shown in Algorithm 10. Each  $\text{AlgBB}^{(i)}$  is an algorithm for a particular bandit Blackwell sequential game (see its general definition in Section 5.1). In order to be able to obtain the feedback required in this bandit Blackwell sequential game, we develop our reduction as if we wanted to maximize  $\widehat{f}_{\delta,t}$  instead of  $f_t$  over feasible region  $\mathcal{P}'$  (not  $\mathcal{P}$ ). Since  $\delta$  is very small, this does not impact the regret of the designed algorithm. Specifically, we consider a bandit Blackwell sequential game with  $\mathcal{X} = \Theta = \mathcal{P}'$ ,  $\mathcal{Y} = [-U_{\delta,n}, U_{\delta,n}]^n$ , where  $U_{\delta,n} = \frac{n}{\delta}$  (i.e., an upper-bound on the absolute value of each coordinate of the vector  $\nabla \widehat{f}_{\delta,t}(\mathbf{z}^{(i-1)})$ ). We further define the vector payoff function  $\mathbf{p}(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} -\mathbf{x} \cdot \mathbf{y} \\ \mathbf{y} \end{bmatrix}$  and the target set  $S = \text{Cone}(\mathbf{1} \oplus \mathcal{P})^\circ$ . In this reduction, we use a different AdvB (compared to the one used for the full-information setting). We define  $\text{AdvB}(\mathbf{z}^{(i-1)}, f)$  to be  $\nabla \widehat{f}_{\delta,t}(\mathbf{z}^{(i-1)})$ , and hence we have  $\mathbf{y}_t^{(i)} = \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t) = \nabla \widehat{f}_{\delta,t}(\mathbf{z}_t^{(i-1)})$ . Moreover, in each round, we have  $\mathbf{x}_t^{(i)} = \boldsymbol{\rho}_t^{(i)}$ .

As stated earlier, at the core of our bandit algorithm, we need to provide the right feedback when subproblem  $i$  explores in round  $t$ . Such a feedback should be an unbiased estimator of

$$\mathbf{p} \left( \boldsymbol{\rho}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t) \right) = \begin{bmatrix} -\boldsymbol{\rho}_t^{(i)} \cdot \nabla \widehat{f}_{\delta,t}(\mathbf{z}_t^{(i-1)}) \\ \nabla \widehat{f}_{\delta,t}(\mathbf{z}_t^{(i-1)}) \end{bmatrix}.$$

Here,  $\boldsymbol{\rho}_t^{(i)}$  is the output of  $\text{AlgBB}^{(i)}$ . Denoting the desired unbiased estimator for subproblem  $i$  in round  $t$  by  $\widehat{\boldsymbol{\rho}}_t^{(i)}$ , we construct  $\widehat{\boldsymbol{\rho}}_t^{(i)}$  with the help of our exploration sampling device, denoted by  $\text{ExpS}$ . Given the

current point  $\mathbf{z}_t^{(i-1)} \in \mathcal{P}'$  received from the previous sub-problem and  $\boldsymbol{\rho}_t^{(i)}$  as the current decision of AlgBB for exploration, ExpS returns  $(\mathbf{w}_{\text{exp}}, \mathbf{z}_{\text{exp}})$  such that (i)  $\mathbf{z}_{\text{exp}} \in \mathcal{P}$ , i.e., it is feasible, and (ii)  $\hat{\mathbf{p}}_t^{(i)} = f_t(\mathbf{z}_{\text{exp}})\mathbf{w}_{\text{exp}}$  is a desired unbiased estimator, i.e., we have

$$\mathbb{E}\left[\hat{\mathbf{p}}_t^{(i)}\right] = \mathbb{E}[f_t(\mathbf{z}_{\text{exp}})\mathbf{w}_{\text{exp}}] = \mathbf{p}\left(\boldsymbol{\rho}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t)\right).$$

Note that the above unbiased estimator can be constructed by only evaluating the function  $f_t$  at a feasible point  $\mathbf{z}_{\text{exp}}$  in  $\mathcal{P}$ . It only remains to show how to construct such an exploration sampling device.

In order to satisfy the above properties (i) and (ii), the exploration sampling device ExpS works as follows: It first uniformly at random picks a point  $\mathbf{u}$  from the sphere  $S_2^{(n-1)}$ . Then it sets

$$\mathbf{w}_{\text{exp}} = \frac{n}{\delta} \begin{bmatrix} -\boldsymbol{\rho}_t^{(i)} \cdot \mathbf{u} \\ \mathbf{u} \end{bmatrix}, \quad \mathbf{z}_{\text{exp}} = \mathbf{z}_t^{(i-1)} + \delta \mathbf{u}.$$

The point  $\mathbf{z}_{\text{exp}}$  is a feasible point in  $\mathcal{P}$ , because  $\mathbf{z}_t^{(i-1)} \in \mathcal{P}'$ . Moreover, we have:

$$\mathbb{E}[f_t(\mathbf{z}_{\text{exp}})\mathbf{w}_{\text{exp}}] = \mathbb{E}\left[\begin{bmatrix} -\boldsymbol{\rho}_t^{(i)} \cdot \frac{n}{\delta} f_t(\mathbf{z}_t^{(i-1)} + \delta \mathbf{u}) \\ \frac{n}{\delta} f_t(\mathbf{z}_t^{(i-1)} + \delta \mathbf{u}) \end{bmatrix} \mathbf{u}\right] = \begin{bmatrix} -\boldsymbol{\rho}_t^{(i)} \cdot \nabla \widehat{f}_{\delta,t}(\mathbf{z}_t^{(i-1)}) \\ \nabla \widehat{f}_{\delta,t}(\mathbf{z}_t^{(i-1)}) \end{bmatrix} = \mathbf{p}\left(\boldsymbol{\rho}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t)\right),$$

which shows  $\hat{\mathbf{p}}_t^{(i)} = f_t(\mathbf{z}_{\text{exp}})\mathbf{w}_{\text{exp}}$  is our desired unbiased estimator.

Given the above bandit Blackwell sequential game, since we have shown that  $S$  is response-satisfiable, there exists a polynomial-time algorithm AlgBB, based on Theorem 3 and its proof in Section D.1, such that for any  $\mathbf{y}_1, \dots, \mathbf{y}_T \in \mathcal{Y}$  and exploration probability  $q \in [0, 1]$ , AlgBB can generate a sequence of decisions  $\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathcal{X}$  such that

$$d_\infty\left(\frac{1}{T} \sum_{t=1}^T \mathbf{p}(\mathbf{x}_t, \mathbf{y}_t), S\right) + \mathbb{E}\left[\frac{1}{T} D(\mathbf{p}) \cdot (\# \text{ explore})\right] \leq O\left(\frac{1}{qT} D(\hat{\mathbf{p}}) (\log d_{\mathbf{p}})^{1/2} (Tq)^{1/2}\right) + O(D(\mathbf{p})q), \quad (30)$$

as shown in the proof of Theorem 3 in Section D.1. Here, the diameter of  $\hat{\mathbf{p}}$ , i.e.,  $D(\hat{\mathbf{p}})$ , is  $\frac{n}{\delta} \max\{1, R_{\mathcal{P}}\}$ . This is because for any  $\mathbf{z} \in \mathcal{P}$ ,  $\mathbf{u} \in S_2^{(n-1)}$ ,  $\left\|\left(\frac{n}{\delta} f(\mathbf{z})\mathbf{u}\right)\right\|_{+\infty} \leq \frac{n}{\delta} \cdot 1 \cdot 1 = \frac{n}{\delta}$  and  $|\boldsymbol{\rho} \cdot \frac{n}{\delta} f(\mathbf{z} + \delta \mathbf{u})\mathbf{u}| \leq \|\boldsymbol{\rho}\|_2 \left\|\frac{n}{\delta} f(\mathbf{z} + \delta \mathbf{u})\mathbf{u}\right\|_2 \leq R_{\mathcal{P}} \frac{n}{\delta}$ , and hence  $D(\hat{\mathbf{p}}) = \frac{n}{\delta} \max\{1, R_{\mathcal{P}}\}$ . We would like to highlight that based on our definition of the actual vector payoff, here, we similarly have  $D(\mathbf{p}) = O\left(R_{\mathcal{P}} \frac{n}{\delta}\right)$ .

We summarize the result for the bandit setting in the following theorem.

**THEOREM 10. (Bandit learning for strong-DR monotone SM maximization in a downward closed convex set)** Let  $\mathcal{P}$  be a downward closed convex polytope with  $\ell_2$  diameter  $R_{\mathcal{P}}$ . Consider the online learning problem of maximizing strong-DR monotone submodular functions on  $\mathcal{P}$ , which are  $L$ -Lipschitz smooth and have bounded gradient of  $U$  for any point  $\mathbf{z} \in \mathcal{P}$ . Then, under bandit setting, with the total number of iterations  $N = O\left(T^{1/6} R_{\mathcal{P}} L^{1/2} (\log(n))^{-1/6}\right)$ , the explore probability  $q = O\left(T^{-1/3} (\log(n))^{1/3}\right)$ , and  $\delta = O\left(T^{-1/6} R_{\mathcal{P}}^{1/2} n^{1/2} (\log(n)^{1/6}) U^{-1/2}\right)$ , Algorithm 10 with  $N$  iterations obtains  $O\left(T^{5/6} (\log(n))^{1/6} R_{\mathcal{P}} \left(n R_{\mathcal{P}}^{1/2} U^{1/2} + L^{1/2}\right)\right) \gamma_N$ -regret, where  $\gamma_N = 1 - (1 - 1/N)^N$ .

Notice that our bandit regret is  $O(T^{5/6})$  in terms of  $T$ , which improves the previous regret bound in Zhang et al. (2020),  $O(T^{8/9})$ , in terms of  $T$ . We now present the proof of the theorem.

**Algorithm 10:** Bandit online learning algorithm for Algorithm 8

**Input:** A sequence of  $L$ -Lipschitz smooth strong-DR monotone submodular functions  $f_1, \dots, f_T : \mathcal{P} \rightarrow [0, 1]$  such that  $\|\nabla f_t(\mathbf{z})\|_\infty \leq U$  for all  $\mathbf{z} \in \mathcal{P}, t \in [T]$ ; parameter  $\delta$ ;  $N$  bandit Blackwell algorithms  $\{\text{AlgBB}^{(i)}\}_{i=1, \dots, N}$ , where each  $\text{AlgBB}^{(i)}$  is an online algorithm for a bandit Blackwell sequential game with  $\mathcal{X} = \mathcal{P}'$  (where  $\mathcal{P}' = (1 - \alpha)\mathcal{P} + \delta \mathbf{1}$ ,  $\alpha = \frac{\sqrt{n+1}}{r}\delta$ , and  $r$  is the largest positive real number such that  $r(B_2^n \cap \mathbb{R}_{\geq 0}^n) \subseteq \mathcal{P}$ ),  $\mathcal{Y} = [-U_{\delta, n}, U_{\delta, n}]^n$  with  $U_{\delta, n} = \frac{n}{\delta}$ ,  $\mathbf{p}(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} -\mathbf{x} \cdot \mathbf{y} \\ \mathbf{y} \end{bmatrix}$ , and  $S = \text{Cone}(\mathbf{1} \oplus \mathcal{P})^\circ$ . Whenever called,  $\text{AlgBB}^{(i)}$  explores with probability  $q$ .

**Output:**  $\mathbf{z}_1, \dots, \mathbf{z}_T \in \mathcal{P}$

Initialize  $N$  parallel instances  $\{\text{AlgBB}^{(i)}\}_{i=1}^N$ .

**for**  $t = 1, \dots, T$  **do**

Initialize  $\mathbf{z}_t^{(0)} \leftarrow \mathbf{0}$ .

**for** iteration  $i = 1, 2, \dots, N$  **do**

Choose the direction  $\boldsymbol{\rho}_t^{(i)} \in \mathcal{P}'$  and the exploration signal  $\pi_t^{(i)} \in \{\text{YES}, \text{NO}\}$  by querying  $\text{AlgBB}^{(i)}$  given the update parameters and vector payoffs  $\hat{\mathbf{p}}_t^{(i)}$  of exploration rounds prior to round  $t$  in the bandit Blackwell sequential game of subproblem  $i$ , i.e.  $(\boldsymbol{\rho}_t^{(i)}, \pi_t^{(i)}) \leftarrow \text{AlgBB}^{(i)}(\boldsymbol{\rho}_1^{(i)}, \dots, \boldsymbol{\rho}_{t-1}^{(i)}, \{\hat{\mathbf{p}}_\tau^{(i)}\}_{\tau \leq t-1}, \pi_\tau^{(i)} = \text{YES})$ .

**if**  $\pi_t^{(i)} = \text{YES}$ , **then**

Play the exploration point  $\mathbf{z}_t \leftarrow \mathbf{z}_t^{(i)} + \delta \mathbf{u}$  where  $\mathbf{u}$  is drawn uniformly at random from the unit  $(n-1)$  dimensional sphere  $S^{(n-1)}$ .

*//Bandit information feedback: observe  $f_t(\mathbf{z}_t + \delta \mathbf{u})$ .*

Give payoff vector feedback  $\hat{\mathbf{p}}_t^{(i)} = \begin{bmatrix} -\boldsymbol{\rho}_t^{(i)} \cdot \frac{n}{\delta} f_t(\mathbf{z}_t + \delta \mathbf{u}) \mathbf{u} \\ \frac{n}{\delta} f_t(\mathbf{z}_t + \delta \mathbf{u}) \mathbf{u} \end{bmatrix}$  to  $\text{AlgBB}^{(i)}$ . Skip immediately to the beginning of the next round  $t+1$ .

*//Note that  $\mathbb{E}[\hat{\mathbf{p}}_t^{(i)}] = \mathbf{p}(\mathbf{x}_t^{(i)}, \mathbf{y}_t^{(i)})$ , where  $\mathbf{x}_t^{(i)} = \boldsymbol{\rho}_t^{(i)}$  and*

*$\mathbf{y}_t^{(i)} = \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t) = \nabla \hat{f}_{\delta, t}(\mathbf{z}_t^{(i)})$ . Here,  $\hat{f}_{\delta, t}(\mathbf{z}_t^{(i)}) = \mathbb{E}_{\mathbf{v} \sim \text{Unif}(B_2^{\mathcal{P}'})}[f_t(\mathbf{z}_t^{(i)} + \delta \mathbf{v})]$ .*

Set  $\mathbf{z}_t^{(i)} \leftarrow \mathbf{z}_t^{(i-1)} + \frac{1}{N} \boldsymbol{\rho}_t^{(i)}$ .

Play the final point  $\mathbf{z}_t \leftarrow \mathbf{z}_t^{(N)}$ , receive the bandit feedback  $f_t(\mathbf{z}_t)$  and ignore them.

*Proof of Theorem 10.* Define  $\mathbf{x}^* = \arg \max_{\mathbf{x} \in \mathcal{P}} \sum_{t=1}^T f_t(\mathbf{x})$  and  $\mathbf{x}_\delta^* = \arg \max_{\mathbf{x} \in \mathcal{P}'} \sum_{t=1}^T \hat{f}_{\delta, t}(\mathbf{x})$ . Since we are running AlgBB as if it were for the function  $\hat{f}_{\delta, t}$  over the feasible region  $\mathcal{P}'$ , our total  $\gamma_N$ -regret can be

written as

$$\begin{aligned}
\mathbb{E} \left[ \sum_{t=1}^T \gamma_N f_t(\mathbf{x}^*) - f_t(\mathbf{z}_t) \right] &= \mathbb{E} \left[ \sum_{t=1}^T \gamma_N (f_t(\mathbf{x}^*) - f_t(\mathbf{x}_\delta^*) + f_t(\mathbf{x}_\delta^*)) - \sum_{t=1}^T f_t(\mathbf{z}_t) \right] \\
&= \mathbb{E} \left[ \gamma_N \sum_{t=1}^T (f_t(\mathbf{x}^*) - f_t(\mathbf{x}_\delta^*)) + \gamma_N \sum_{t=1}^T (f_t(\mathbf{x}_\delta^*) - \widehat{f}_{\delta,t}(\mathbf{x}_\delta^*)) + \sum_{t=1}^T (\gamma_N \widehat{f}_{\delta,t}(\mathbf{x}_\delta^*) - \widehat{f}_{\delta,t}(\mathbf{z}_t)) \right. \\
&\quad \left. + \sum_{t=1}^T (\widehat{f}_{\delta,t}(\mathbf{z}_t) - f_t(\mathbf{z}_t)) \right] \\
&\stackrel{(1)}{\leq} \mathbb{E} \left[ \gamma_N \sum_{t=1}^T (f_t(\mathbf{x}^*) - f_t(\mathbf{x}_\delta^*)) + \gamma_N TU\delta + \sum_{t=1}^T (\gamma_N \widehat{f}_{\delta,t}(\mathbf{x}_\delta^*) - \widehat{f}_{\delta,t}(\mathbf{z}_t)) + TU\delta \right], \tag{31}
\end{aligned}$$

where Inequality (1) holds because we have  $|f_t(\mathbf{z}) - \widehat{f}_{\delta,t}(\mathbf{z})| \leq U\delta$  for any  $\mathbf{z} \in \mathcal{P}$ . We are now left with bounding the first and third terms in Equation (31).

To bound the first term in Equation (31), let  $\mathbf{x}' \in \mathcal{P}'$  be the point in  $\mathcal{P}'$  that is closest to  $\mathbf{x}^*$ , i.e.,  $\|\mathbf{x}' - \mathbf{x}^*\|_2 = d(\mathbf{x}', \mathbf{x}^*) = d(\mathbf{x}^*, \mathcal{P}') \leq d(\mathcal{P}, \mathcal{P}')$ , where  $d(\mathbf{x}, \mathcal{P}') = \min_{\mathbf{v} \in \mathcal{P}'} \|\mathbf{x} - \mathbf{v}\|_2$  and  $d(\mathcal{P}, \mathcal{P}') = \max_{\mathbf{v} \in \mathcal{P}} d(\mathbf{x}, \mathcal{P}')$ .

We then have

$$\begin{aligned}
\sum_{t=1}^T (f_t(\mathbf{x}^*) - f_t(\mathbf{x}_\delta^*)) &= \sum_{t=1}^T (f_t(\mathbf{x}^*) - f_t(\mathbf{x}')) + \sum_{t=1}^T (f_t(\mathbf{x}') - f_t(\mathbf{x}_\delta^*)) \\
&\stackrel{(1)}{\leq} U \sum_{t=1}^T \|\mathbf{x}^* - \mathbf{x}'\|_2 + 0 \\
&\leq UTd(\mathcal{P}, \mathcal{P}') \\
&\stackrel{(2)}{\leq} UT \left( \sqrt{n} \left( \frac{R_{\mathcal{P}}}{r} + 1 \right) + \frac{R_{\mathcal{P}}}{r} \right) \delta, \tag{32}
\end{aligned}$$

where Inequality (1) holds because  $f_t$  is  $U$ -Lipschitz continuous and  $\mathbf{x}_\delta^*$  maximizes  $f_t$  in  $\mathcal{P}'$ , while Inequality (2) holds because

$$d(\mathcal{P}, \mathcal{P}') \leq \left( \sqrt{n} \left( \frac{R_{\mathcal{P}}}{r} + 1 \right) + \frac{R_{\mathcal{P}}}{r} \right) \delta$$

from Lemma 1 in Zhang et al. (2020).

To bound the third term of Equation (31), i.e.,  $\sum_{t=1}^T (\gamma_N \widehat{f}_{\delta,t}(\mathbf{x}_\delta^*) - \widehat{f}_{\delta,t}(\mathbf{z}_t))$ , we follow the same lines as in the proof for Theorem 4, by adapting it as if the underlying functions are  $\widehat{f}_{\delta,t}$ . We then employ the result in Equation (26) for a new value of  $h(T)$  to show that Algorithm 10 works. Let  $\mathcal{T}_i \subseteq [T]$  be the set of rounds where  $\text{AlgBB}^{(i)}$  is invoked. Note that  $\mathcal{T}_i$  is a random set since it depends on the realization of binary signals  $\{\pi_t^{(i')}\}_{i' \in [i-1], t \in [T]}$ . Consider a subproblem  $i \in [N]$ . Fix a particular realization of set  $\mathcal{T}_i$ . By using the upper-bound in Equation (30) and the fact that  $D(\hat{\mathbf{p}}) = O\left(\frac{n}{\delta} R_{\mathcal{P}}\right)$  and  $D(\mathbf{p}) = O\left(\frac{n}{\delta} R_{\mathcal{P}}\right)$ , we have

$$\begin{aligned}
d_\infty \left( \frac{1}{|\mathcal{T}_i|} \sum_{t \in \mathcal{T}_i} \mathbf{p} \left( \boldsymbol{\rho}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t) \right), S \right) &\leq O \left( \frac{1}{q|\mathcal{T}_i|} D(\hat{\mathbf{p}}) \log(d_{\mathbf{p}})^{1/2} (|\mathcal{T}_i|q)^{1/2} \right) + O(D(\mathbf{p})q) \\
&= O \left( \frac{1}{q|\mathcal{T}_i|} \frac{n}{\delta} R_{\mathcal{P}} (\log n)^{1/2} (|\mathcal{T}_i|q)^{1/2} \right) + O \left( \frac{n}{\delta} R_{\mathcal{P}} q \right).
\end{aligned}$$

Since  $\mathbf{p} \left( \boldsymbol{\rho}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t) \right) = \text{PAYOFF} \left( \boldsymbol{\rho}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t \right)$  by definition, we have

$$O \left( \frac{1}{q|\mathcal{T}_i|} \frac{n}{\delta} R_{\mathcal{P}} (\log n)^{1/2} (|\mathcal{T}_i|q)^{1/2} \right) + O \left( \frac{n}{\delta} q \right) \geq d_\infty \left( \frac{1}{|\mathcal{T}_i|} \sum_{t \in \mathcal{T}_i} \mathbf{p} \left( \boldsymbol{\rho}_t^{(i)}, \text{AdvB}(\mathbf{z}_t^{(i-1)}, f_t) \right), S \right)$$

$$\begin{aligned}
&= d_\infty \left( \frac{1}{|\mathcal{T}_i|} \sum_{t \in \mathcal{T}_i} \text{PAYOFF} \left( \boldsymbol{\rho}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t \right), S \right) \\
&\geq \frac{1}{\sqrt{n+1} \sqrt{1+R_{\mathcal{P}}^2}} \frac{1}{|\mathcal{T}_i|} \left( \max_{\boldsymbol{\rho} \in \mathcal{P}'} \sum_{t \in \mathcal{T}_i} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho} - \sum_{t \in \mathcal{T}_i} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}_t^{(i)} \right),
\end{aligned} \tag{33}$$

where the last line follows from the equivalent of Equation (27) and Equation (28), which is still true because by shrinking the domain for  $\{\mathbf{z}_t^{(i)}\}_{1 \leq i \leq N, 1 \leq t \leq T}$  to be  $\mathcal{P}'$ , all of the chosen points are guaranteed to be in  $\mathcal{P}$ , so  $\mathcal{P}' \subseteq \mathcal{P}$  and the diameter of  $\mathcal{P}'$  is upper-bounded by the diameter of  $\mathcal{P}$ ,  $R_{\mathcal{P}}$  as  $\mathcal{P}' \subseteq \mathcal{P}$ .

Let  $\mathcal{T}^-$  be the set of rounds where no AlgBB<sup>(i)</sup> explored and  $\mathcal{T}^+$  be the set of rounds where some AlgBB<sup>(i)</sup> explored. See that  $\mathcal{T}^- \subseteq \mathcal{T}_i$  because if no algorithm explores, AlgBB<sup>(i)</sup> will be invoked. For a set of rounds  $\mathcal{T}$ , let  $\boldsymbol{\rho}^*(\mathcal{T}) = \arg \max_{\boldsymbol{\rho} \in \mathcal{P}'} \sum_{t \in \mathcal{T}} \boldsymbol{\rho} \cdot \nabla \widehat{f}_{\delta,t}(\mathbf{z}_t^{(i-1)})$ . Also, let  $M_i$  be the number of rounds that AlgBB<sup>(i)</sup> explores out of the  $|\mathcal{T}_i|$  rounds it is invoked. For subproblem  $i$ , we have

$$\begin{aligned}
&\mathbb{E} \left[ \left( \max_{\boldsymbol{\rho} \in \mathcal{P}'} \sum_{t \in \mathcal{T}^-} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho} - \sum_{t \in \mathcal{T}^-} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}_t^{(i)} \right) \middle| \mathcal{T}_i \right] \\
&= \mathbb{E} \left[ \left( \sum_{t \in \mathcal{T}^-} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}^*(\mathcal{T}^-) - \sum_{t \in \mathcal{T}^-} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}_t^{(i)} \right) \middle| \mathcal{T}_i \right] \\
&= \mathbb{E} \left[ \left( \sum_{t \in \mathcal{T}_i} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}^*(\mathcal{T}^-) - \sum_{t \in \mathcal{T}_i} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}_t^{(i)} \right) \middle| \mathcal{T}_i \right] - \mathbb{E} \left[ \left( \sum_{t \in \mathcal{T}_i \setminus \mathcal{T}^-} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}^*(\mathcal{T}^-) - \sum_{t \in \mathcal{T}_i \setminus \mathcal{T}^-} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}_t^{(i)} \right) \middle| \mathcal{T}_i \right] \\
&\leq \mathbb{E} \left[ \left( \sum_{t \in \mathcal{T}_i} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}^*(\mathcal{T}_i) - \sum_{t \in \mathcal{T}_i} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}_t^{(i)} \right) \middle| \mathcal{T}_i \right] - \mathbb{E} \left[ \left( \sum_{t \in \mathcal{T}_i \setminus \mathcal{T}^-} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}^*(\mathcal{T}^-) - \sum_{t \in \mathcal{T}_i \setminus \mathcal{T}^-} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}_t^{(i)} \right) \middle| \mathcal{T}_i \right] \\
&\stackrel{(1)}{\leq} O \left( \frac{R_{\mathcal{P}} \sqrt{n}}{q} \frac{n}{\delta} R_{\mathcal{P}} (\log n)^{1/2} (|\mathcal{T}_i| q)^{1/2} \right) + O \left( R_{\mathcal{P}} \sqrt{n} |\mathcal{T}_i| \frac{n}{\delta} R_{\mathcal{P}} q \right) + D(\mathbf{p}) \mathbb{E}[M_i | \mathcal{T}_i] \\
&\stackrel{(2)}{\leq} O \left( \frac{R_{\mathcal{P}}^2 n^{3/2}}{q^{1/2} \delta} (\log n)^{1/2} |\mathcal{T}_i|^{1/2} \right) + O \left( \frac{R_{\mathcal{P}}^2 n^{3/2}}{\delta} |\mathcal{T}_i| q \right) + O \left( \frac{n}{\delta} R_{\mathcal{P}} q |\mathcal{T}_i| \right) \\
&\stackrel{(3)}{\leq} O \left( \frac{R_{\mathcal{P}}^2 n^{3/2}}{q^{1/2} \delta} (\log n)^{1/2} T^{1/2} \right) + O \left( \frac{R_{\mathcal{P}}^2 n^{3/2}}{\delta} T q \right) + O \left( \frac{n}{\delta} R_{\mathcal{P}} q T \right),
\end{aligned}$$

where all of the expectations are with respect to  $\mathbf{z}_t^{(i-1)}, t \in \mathcal{T}_i$ . Here, Inequality (1) follows from Equation (33) and the fact that  $|\mathcal{T}_i \setminus \mathcal{T}^-| \leq M_i$ . Moreover, for any  $t \in [T]$ ,  $\boldsymbol{\rho}_1, \boldsymbol{\rho}_2 \in \mathcal{P} \subseteq [0, 1]^n$ ,  $\mathbf{y}_t^{(i)} \cdot (\boldsymbol{\rho}_1 - \boldsymbol{\rho}_2) \leq D_\infty(\mathbf{y}_t^{(i)}) \leq U$  since  $\mathcal{P}$  is closed-down. Inequality (2) holds because  $\mathbb{E}[M_i | \mathcal{T}_i] = q |\mathcal{T}_i|$ . Inequality (3) holds because  $|\mathcal{T}_i| \leq T$ . Fixing a function  $h(\cdot)$ , see that given a fixed  $\mathcal{T}_i$ ,

$$\mathbb{E} \left[ \left( \max_{\boldsymbol{\rho} \in \mathcal{P}'} \sum_{t \in \mathcal{T}^-} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho} - \sum_{t \in \mathcal{T}^-} \mathbf{y}_t^{(i)} \cdot \boldsymbol{\rho}_t^{(i)} \right) \middle| \mathcal{T}_i \right] \leq h(|\mathcal{T}^-|)$$

is equivalent to

$$\begin{aligned}
&\mathbb{E} \left[ \left( \max_{\boldsymbol{\rho} \in \mathcal{P}'} \sum_{t \in \mathcal{T}^-} \nabla \widehat{f}_{\delta,t}(\mathbf{z}_t^{(i-1)}) \cdot \boldsymbol{\rho} - \sum_{t \in \mathcal{T}^-} \nabla \widehat{f}_{\delta,t}(\mathbf{z}_t^{(i-1)}) \cdot \boldsymbol{\rho}_t^{(i)} \right) \right] \leq h(|\mathcal{T}^-|) \\
&\Leftrightarrow \mathbb{E} \left[ \sum_{t \in \mathcal{T}^-} \nabla \widehat{f}_{\delta,t}(\mathbf{z}_t^{(i-1)}) \cdot \boldsymbol{\rho}_t^{(i)} \right] \geq \mathbb{E} \left[ \max_{\boldsymbol{\rho} \in \mathcal{P}'} \sum_{t \in \mathcal{T}^-} \nabla \widehat{f}_{\delta,t}(\mathbf{z}_t^{(i-1)}) \cdot \boldsymbol{\rho} \right] - h(|\mathcal{T}^-|).
\end{aligned}$$

We can then directly invoke the result in Equation (26) for a fixed  $\mathcal{T}^-$ , with  $h(|\mathcal{T}^-|)$  being replaced with its upper-bound above, i.e.,  $O\left(\frac{R_{\mathcal{P}}^2 n^{3/2}}{q^{1/2}\delta}(\log n)^{1/2}T^{1/2} + \frac{R_{\mathcal{P}}^2 n^{3/2}}{\delta}Tq + \frac{n}{\delta}R_{\mathcal{P}}qT\right)$ . Then for any  $\mathbf{z}_\delta^* \in \mathcal{P}'$ , we have

$$\begin{aligned} \mathbb{E}\left[\sum_{t \in \mathcal{T}^-} \widehat{f}_{\delta,t}(\mathbf{z}_t^{(N)})\right] &\geq \gamma_N \mathbb{E}\left[\sum_{t \in \mathcal{T}^-} \widehat{f}_{\delta,t}(\mathbf{z}_\delta^*)\right] - O\left(\frac{R_{\mathcal{P}}^2 n^{3/2}}{q^{1/2}\delta}(\log n)^{1/2}T^{1/2}\right) - O\left(\frac{R_{\mathcal{P}}^2 n^{3/2}}{\delta}Tq\right) \\ &\quad - O\left(\frac{n}{\delta}R_{\mathcal{P}}qT\right) - \frac{|\mathcal{T}^-|LR_{\mathcal{P}}^2}{2N} \\ &\geq \gamma_N \mathbb{E}\left[\sum_{t \in \mathcal{T}^-} \widehat{f}_{\delta,t}(\mathbf{z}_\delta^*)\right] - O\left(\frac{R_{\mathcal{P}}^2 n^{3/2}}{q^{1/2}\delta}(\log n)^{1/2}T^{1/2}\right) - O\left(\frac{R_{\mathcal{P}}^2 n^{3/2}}{\delta}Tq\right) \\ &\quad - O\left(\frac{n}{\delta}R_{\mathcal{P}}qT\right) - O\left(\frac{TLR_{\mathcal{P}}^2}{N}\right), \end{aligned} \quad (34)$$

We further show that Algorithm 10 does not explore too often among its subproblems, specifically

$$\mathbb{E}[|\mathcal{T}^+|] = \mathbb{E}\left[\sum_{i=1}^N M_i\right] \leq \sum_{i=1}^N O(q\mathbb{E}[T_i]) \leq \sum_{i=1}^N O(qT),$$

and since functions  $\widehat{f}_{\delta,t}$  output values that are in  $[0, 1]$ , for the remaining rounds  $\mathcal{T}^+$  we have:

$$\mathbb{E}\left[\sum_{t \in \mathcal{T}^+} \widehat{f}_{\delta,t}(\mathbf{z}_t)\right] \geq \gamma_N \cdot \mathbb{E}\left[\sum_{t \in \mathcal{T}^+} \widehat{f}_{\delta,t}(\mathbf{z}_\delta^*)\right] - O(NqT). \quad (35)$$

Summing Equation (34) and Equation (35), we have

$$\begin{aligned} \mathbb{E}\left[\gamma_N \sum_{t \in [T]} \widehat{f}_{\delta,t}(\mathbf{z}_\delta^*) - \sum_{t \in [T]} \widehat{f}_{\delta,t}(\mathbf{z}_t)\right] &\leq O\left(\frac{R_{\mathcal{P}}^2 n^{3/2}}{q^{1/2}\delta}(\log n)^{1/2}T^{1/2}\right) + O\left(\frac{R_{\mathcal{P}}^2 n^{3/2}}{\delta}Tq\right) \\ &\quad + O\left(\frac{n}{\delta}R_{\mathcal{P}}qT\right) + O\left(\frac{TLR_{\mathcal{P}}^2}{N}\right) + O(NqT), \end{aligned} \quad (36)$$

which we can use to bound the third term in Equation (31). Revisiting the total regret, and substituting Equation (32) and Equation (36) to our total  $\gamma_N$ -regret in Equation (31), we then have

$$\begin{aligned} \mathbb{E}\left[\sum_{t=1}^T \gamma_N f_t(\mathbf{x}^*) - f_t(\mathbf{z}_t)\right] &\leq O\left(TU\delta\sqrt{n}R_{\mathcal{P}} + \frac{R_{\mathcal{P}}^2 n^{3/2}}{q^{1/2}\delta}(\log n)^{1/2}T^{1/2} + \frac{R_{\mathcal{P}}^2 n^{3/2}}{\delta}Tq\right. \\ &\quad \left. + \frac{TLR_{\mathcal{P}}^2}{N} + NqT\right). \end{aligned}$$

We finally tune parameters  $q$ ,  $\delta$ , and  $N$  for the best regret bound. Specifically, by substituting  $q = O\left(T^{-1/3}(\log(n))^{1/3}\right)$ ,  $\delta = O\left(T^{-1/6}R_{\mathcal{P}}^{1/2}n^{1/2}(\log(n))^{1/6}U^{-1/2}\right)$ , and  $N = O\left(T^{1/6}R_{\mathcal{P}}L^{1/2}(\log(n))^{-1/6}\right)$ , we get the total regret of

$$O\left(T^{5/6}(\log(n))^{1/6}R_{\mathcal{P}}\left(nR_{\mathcal{P}}^{1/2}U^{1/2} + L^{1/2}\right)\right).$$

This completes the proof. ■

REMARK 6. As we showed in Theorem 10, we can obtain sub-linear approximate regret with an optimal approximation factor for the general strong-DR Lipschitz continuous smooth submodular functions under the bandit feedback. As multi-linear extension of a set submodular function satisfies these conditions, it seems possible to extend our result in Section F.2 to this setting. However, there is a technical caveat and

that is generating an unbiased estimator for  $f_t^{\text{MLE}}(\mathbf{x})$  given query access at *feasible* sets to the set function  $f_t$ , which is needed to obtain the required bandit feedback for Algorithm 10 when it is run on the sequence of functions  $f_1^{\text{MLE}}, \dots, f_T^{\text{MLE}}$ . Given a point  $\mathbf{x} \in \mathcal{P}$ , where  $\mathcal{P}$  is the matroid polytope, the pipage rounding algorithm will sample a randomized set  $S$  such that  $\mathbb{E}[f_t(S)] \geq f_t^{\text{MLE}}(\mathbf{x})$ . To have the equality, and hence an unbiased estimator for  $f_t^{\text{MLE}}(\mathbf{x})$ , one possible way is placing each element independently with probability  $x_i$  in random set  $S$ ; but then this set is not necessarily an independent set and hence we are not allowed to return it as the algorithm in our problem. We pose whether our approach can be extended to the bandit setting as an open problem.

## Appendix G: Proofs and Remarks of Section 6.1 – Product Ranking and Sequential Submodular Maximization

In this appendix, we give the missing proofs of the results from Section 6.1.

### G.1. Proof of Theorem 5

*Proof.* We show that our meta Algorithm 4 works by verifying the following conditions.

(i) *Algorithm 5 is an extended  $(\frac{1}{2}, \frac{1}{2})$ -robust approximation algorithm.* We need to show that if the following equation holds for some function  $h$ :

$$\forall j \in [n], \quad \left[ \sum_{t=1}^T \text{PAYOFF}(\tilde{\theta}_t^{(i)}, \boldsymbol{\pi}_t^{(i-1)}, f_t) \right]_j \geq -h(T),$$

then we must have

$$\forall \boldsymbol{\pi}^* \in \Pi, \quad \sum_{t=1}^T \mathbb{E}[f_t(\boldsymbol{\pi}_t)] \geq \frac{1}{2} \sum_{t=1}^T f_t(\boldsymbol{\pi}^*) - \frac{1}{2} nh(T). \quad (37)$$

Recall that for any  $j \in [n]$ , we have  $[\text{PAYOFF}(\boldsymbol{\theta}^{(i)}, \boldsymbol{\pi}^{(i-1)}, f)]_j = \boldsymbol{\theta}^T \mathbf{y}^{(i)} - [\mathbf{y}^{(i)}]_j$ , where  $\mathbf{y}^{(i)} \triangleq [f(\boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_i) - f(\boldsymbol{\pi}^{(i-1)})]_{j \in [n]}$ . First, we prove several inequalities that will later be used to prove inequality (37). Since each function  $f_{t,i}$  is monotone submodular, we have:

$$\begin{aligned} & \lambda_i \sum_{j=1}^i \left( f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_j\}) - f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_{j-1}\}) \right) \\ & \quad + \lambda_i \sum_{j=1}^i \left( f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_{j-1}, [\boldsymbol{\pi}^*]_j\}) - f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_{j-1}\}) \right) \\ & \stackrel{(1)}{\geq} \lambda_i \sum_{j=1}^i \left( f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_j, [\boldsymbol{\pi}^*]_1, \dots, [\boldsymbol{\pi}^*]_j\}) - f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_{j-1}, [\boldsymbol{\pi}^*]_1, \dots, [\boldsymbol{\pi}^*]_{j-1}\}) \right) \\ & = \lambda_i f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_i, [\boldsymbol{\pi}^*]_1, \dots, [\boldsymbol{\pi}^*]_i\}) - \lambda_i f_{t,i}(\emptyset) \\ & \stackrel{(2)}{\geq} \lambda_i f_{t,i}(\{[\boldsymbol{\pi}^*]_1, \dots, [\boldsymbol{\pi}^*]_i\}), \end{aligned}$$

where inequality (1) follows from submodularity and inequality (2) follows from monotonicity and non-negativity of  $f_{t,i}$ . To be more clear, inequality (1) holds because for each  $j = 1, 2, \dots, i$ , the sum of the marginal values of adding  $[\boldsymbol{\pi}_t]_j$  to  $\boldsymbol{\pi}^{(j-1)}$  and adding  $[\boldsymbol{\pi}^*]_j$  to  $\boldsymbol{\pi}^{(j-1)}$  is greater than equal to the marginal value of adding  $\{[\boldsymbol{\pi}_t]_j, [\boldsymbol{\pi}^*]_j\}$  to  $\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_i, [\boldsymbol{\pi}^*]_1, \dots, [\boldsymbol{\pi}^*]_i\}$  as  $f_{t,i}$  is submodular. Recall that

$$f_t(\boldsymbol{\pi}) \triangleq \lambda_1 f_{t,1}(\{[\boldsymbol{\pi}]_1\}) + \lambda_2 f_{t,2}(\{[\boldsymbol{\pi}]_1, [\boldsymbol{\pi}]_2\}) + \dots + \lambda_n f_{t,n}(\{[\boldsymbol{\pi}]_1, \dots, [\boldsymbol{\pi}]_n\}) \quad \forall t \in [T],$$

so summing the inequalities above for  $i = 1, 2, \dots, n$ , we get:

$$\begin{aligned}
& \sum_{i=1}^n \lambda_i \sum_{j=1}^i \left( f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_j\}) - f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_{j-1}\}) \right) \\
& + \sum_{i=1}^n \lambda_i \sum_{j=1}^i \left( f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_{j-1}, [\boldsymbol{\pi}^*]_j\}) - f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_{j-1}\}) \right) \geq \sum_{i=1}^n \lambda_i f_{t,i}(\{[\boldsymbol{\pi}^*]_1, \dots, [\boldsymbol{\pi}^*]_i\}) \\
& \Leftrightarrow \sum_{j=1}^n \sum_{i=j}^n (\lambda_i f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_j\}) - \lambda_i f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_{j-1}\})) \\
& + \sum_{j=1}^n \sum_{i=j}^n (\lambda_i f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_{j-1}, [\boldsymbol{\pi}^*]_j\}) - \lambda_i f_{t,i}(\{[\boldsymbol{\pi}_t]_1, \dots, [\boldsymbol{\pi}_t]_{j-1}\})) \geq f_t(\boldsymbol{\pi}^*) \\
& \Leftrightarrow \sum_{j=1}^n \left( f_t(\boldsymbol{\pi}_t^{(j)}) - f_t(\boldsymbol{\pi}_t^{(j-1)}) \right) + \sum_{j=1}^n \left( f_t(\boldsymbol{\pi}_t^{(j-1)} + [\boldsymbol{\pi}^*]_j \mathbf{e}_j) - f_t(\boldsymbol{\pi}_t^{(j-1)}) \right) \geq f_t(\boldsymbol{\pi}^*). \tag{38}
\end{aligned}$$

We get the first equivalence by switching the summations. We now use the inequality (38) to prove the final claim, i.e., the desired inequality (37). We have:

$$\begin{aligned}
\sum_{t=1}^T \mathbb{E}[f_t(\boldsymbol{\pi}_t)] &= \frac{1}{2} \sum_{t=1}^T \sum_{i=1}^n 2 \mathbb{E} \left[ f_t(\boldsymbol{\pi}_t^{(i)}) - f_t(\boldsymbol{\pi}_t^{(i-1)}) \right] \\
&= \frac{1}{2} \sum_{t=1}^T \sum_{i=1}^n \mathbb{E} \left[ f_t(\boldsymbol{\pi}_t^{(i)}) - f_t(\boldsymbol{\pi}_t^{(i-1)}) \right] + \frac{1}{2} \sum_{i=1}^n \sum_{t=1}^T \mathbb{E} \left[ f_t(\boldsymbol{\pi}_t^{(i)}) - f_t(\boldsymbol{\pi}_t^{(i-1)}) \right] \\
&\quad - \frac{1}{2} \sum_{i=1}^n \sum_{t=1}^T \mathbb{E} \left[ f_t(\boldsymbol{\pi}_t^{(i)} + [\boldsymbol{\pi}^*]_i \mathbf{e}_i) - f_t(\boldsymbol{\pi}_t^{(i-1)}) \right] + \frac{1}{2} \sum_{i=1}^n \sum_{t=1}^T \mathbb{E} \left[ f_t(\boldsymbol{\pi}_t^{(i)} + [\boldsymbol{\pi}^*]_i \mathbf{e}_i) - f_t(\boldsymbol{\pi}_t^{(i-1)}) \right] \\
&\stackrel{(1)}{=} \frac{1}{2} \sum_{t=1}^T \sum_{i=1}^n \mathbb{E} \left[ f_t(\boldsymbol{\pi}_t^{(i)}) - f_t(\boldsymbol{\pi}_t^{(i-1)}) \right] + \frac{1}{2} \sum_{i=1}^n \sum_{t=1}^T \left[ \text{PAYOFF}(\tilde{\boldsymbol{\theta}}_t^{(i)}, \boldsymbol{\pi}_t^{(i-1)}, f_t) \right]_{[\boldsymbol{\pi}^*]_i} \\
&\quad + \frac{1}{2} \sum_{i=1}^n \sum_{t=1}^T \mathbb{E} \left[ f_t(\boldsymbol{\pi}_t^{(i-1)} + [\boldsymbol{\pi}^*]_i \mathbf{e}_i) - f_t(\boldsymbol{\pi}_t^{(i-1)}) \right] \\
&\stackrel{(2)}{\geq} \frac{1}{2} \sum_{t=1}^T \sum_{i=1}^n \mathbb{E} \left[ f_t(\boldsymbol{\pi}_t^{(i)}) - f_t(\boldsymbol{\pi}_t^{(i-1)}) \right] - \frac{1}{2} nh(T) + \frac{1}{2} \sum_{t=1}^T \sum_{i=1}^n \mathbb{E} \left[ f_t(\boldsymbol{\pi}_t^{(i)} + [\boldsymbol{\pi}^*]_i \mathbf{e}_i) - f_t(\boldsymbol{\pi}_t^{(i-1)}) \right] \\
&\stackrel{(3)}{\geq} \frac{1}{2} \sum_{t=1}^T f_t(\boldsymbol{\pi}^*) - \frac{1}{2} nh(T) .
\end{aligned}$$

In the above chain of inequalities, equality (1) follows from the definition of PAYOFF, inequality (2) follows from our assumption, and inequality (3) follows from inequality (38). Rearranging the terms will finish the proof.

(ii) *Algorithm 5 is bandit Blackwell reducible.* We verify the following conditions based on Definition 9 to show bandit Blackwell reducibility:

- *Algorithm 5 is Blackwell reducible.* For each subproblem, consider an instance  $(\mathcal{X}, \mathcal{Y}, \mathbf{p})$  of Blackwell where  $\mathcal{X} = \Theta = \Delta([n])$  and  $\mathcal{Y} = [-1, 1]^n$ . Our Blackwell adversary function is the marginal increase in the objective function of placing item on position  $i$ ,  $\text{AdvB}(\boldsymbol{\pi}^{(i-1)}, f) = \left[ f(\boldsymbol{\pi}_t^{(i-1)} + j \mathbf{e}_i) - f(\boldsymbol{\pi}_t^{(i-1)}) \right]_{j \in [n]}$ . The biaffine payoff is  $\mathbf{p}(\boldsymbol{\theta}, \mathbf{y}) = \boldsymbol{\theta}^T \mathbf{y} \mathbf{1}_n - \mathbf{y}$ , where  $\mathbf{1}_n$  is an  $n$ -dimensional all ones vector. The target set  $S$  is the non-negative orthant, and it is response-satisfiable since for every adversary's action  $\mathbf{y} \in \mathcal{Y}$ , the strategy  $\boldsymbol{\theta} = \mathbf{e}_{j^*}$  where  $j^* = \arg\max_{j \in [n]} [\mathbf{y}]_j$  results in  $\mathbf{p}(\boldsymbol{\theta}, \mathbf{y}) \geq 0$ .

- An unbiased estimator for the Blackwell payoff function  $\mathbf{p}$  can be constructed. Specifically, we need to construct an exploration sampling device  $U$  that receives  $(\boldsymbol{\theta}, \boldsymbol{\pi}^{(i-1)})$  in subproblem  $i$  and returns  $(\mathbf{w}_{\text{exp}}, \boldsymbol{\pi}_{\text{exp}})$  such that (i) for all  $f \in \mathcal{F}, \boldsymbol{\theta} \in \Theta, \boldsymbol{\pi}^{(i-1)} \in \mathcal{D}, i \in [n]$ :  $\hat{\mathbf{p}}(\boldsymbol{\theta}, \text{AdvB}(\boldsymbol{\pi}^{(i-1)}, f)) = f(\boldsymbol{\pi}_{\text{exp}})\mathbf{w}_{\text{exp}}$ , where  $(\mathbf{w}_{\text{exp}}, \boldsymbol{\pi}_{\text{exp}}) \sim \text{ExpS}(\boldsymbol{\theta}, \boldsymbol{\pi}^{(i-1)})$ , and (ii)  $\hat{\mathbf{p}}$  is an unbiased estimator for the actual payoff, i.e.,  $\forall \boldsymbol{\theta} \in \Theta, \mathbf{y} \in \mathcal{Y} : \mathbb{E}[\hat{\mathbf{p}}(\boldsymbol{\theta}, \mathbf{y})] = \mathbf{p}(\boldsymbol{\theta}, \mathbf{y})$ . The explore sampling device  $U$  works as follows. Given a point  $\boldsymbol{\pi}^{(i-1)} \in \Pi$  and a parameter  $\boldsymbol{\theta} \in \Theta$ , it draws  $j \sim \text{Uniform}\{1, 2, \dots, n\}$  and returns

$$(\mathbf{w}_{\text{exp}}, \boldsymbol{\pi}_{\text{exp}}) = (n(\boldsymbol{\theta}_j \mathbf{1}_n - \mathbf{e}_j), \boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_j). \quad (39)$$

Now,  $\hat{\mathbf{p}}$  is an unbiased estimator of  $\mathbf{p}$  because

$$\begin{aligned} \mathbb{E}[\hat{\mathbf{p}}(\boldsymbol{\theta}, \mathbf{y})] &= \mathbb{E}[\hat{\mathbf{p}}(\boldsymbol{\theta}, \text{AdvB}(\boldsymbol{\pi}^{(i-1)}, f))] = \mathbb{E}[f(\boldsymbol{\pi}_{\text{exp}})\mathbf{w}_{\text{exp}}] \\ &= \mathbb{E}\left[n\boldsymbol{\theta}_j f(\boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_j)\mathbf{1}_n - f(\boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_j)\mathbf{e}_j\right] \\ &\stackrel{(1)}{=} \mathbf{1}_n \sum_{j=1}^n \boldsymbol{\theta}_j f(\boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_j) - \left[f(\boldsymbol{\pi}^{(i-1)} + \mathbf{1e}_i), f(\boldsymbol{\pi}^{(i-1)} + 2\mathbf{e}_i), \dots, f(\boldsymbol{\pi}^{(i-1)} + n\mathbf{e}_i)\right]^T \\ &\stackrel{(2)}{=} \mathbf{1}_n \sum_{j=1}^n \boldsymbol{\theta}_j (f(\boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_j) - f(\boldsymbol{\pi}^{(i-1)})) \\ &\quad - \left[f(\boldsymbol{\pi}^{(i-1)} + \mathbf{1e}_i), f(\boldsymbol{\pi}^{(i-1)} + 2\mathbf{e}_i), \dots, f(\boldsymbol{\pi}^{(i-1)} + n\mathbf{e}_i)\right]^T + f(\boldsymbol{\pi}^{(i-1)})\mathbf{1}_n \\ &= \mathbf{1}_n \sum_{j=1}^n \boldsymbol{\theta}_j (f(\boldsymbol{\pi}^{(i-1)} + j\mathbf{e}_j) - f(\boldsymbol{\pi}^{(i-1)})) \\ &\quad - \left[f(\boldsymbol{\pi}^{(i-1)} + \mathbf{1e}_i) - f(\boldsymbol{\pi}^{(i-1)}), \dots, f(\boldsymbol{\pi}^{(i-1)} + n\mathbf{e}_i) - f(\boldsymbol{\pi}^{(i-1)})\right]^T \\ &= \boldsymbol{\theta}^T \mathbf{y} \mathbf{1}_n - \mathbf{y} = \mathbf{p}(\boldsymbol{\theta}, \text{AdvB}(\boldsymbol{\pi}^{(i-1)}, f)), \end{aligned}$$

where  $\mathbf{y} \triangleq [f(\boldsymbol{\pi}^{(i-1)}(1), \dots, \boldsymbol{\pi}^{(i-1)}(i-1), j) - f(\boldsymbol{\pi}^{(i-1)})]_{j \in [n]}$ .

Here, equation (1) holds because we take  $j \sim \text{Uniform}\{1, 2, \dots, n\}$ , and equation (2) holds because  $\sum_{j=1}^n \boldsymbol{\theta}_j = 1$ . Intuitively, at every round,  $U$  randomly picks one of the items  $j \in [n]$ , and evaluate the marginal benefit of putting element  $j$  on the  $i^{\text{th}}$  position of  $\boldsymbol{\pi}^{(i-1)}$ .

Putting (i) and (ii) altogether, Algorithm 5 is a  $(\frac{1}{2}, \frac{1}{2})$ -extended robust approximation algorithm with  $n$  subproblems. Its payoff diameter  $D(\mathbf{p})$  is  $O(1)$  and its payoff estimator diameter  $D(\hat{\mathbf{p}})$  is  $O(n)$ . The dimension of vector payoffs is also  $d_{\text{payoff}} = n$ . It is also bandit Blackwell reducible, hence from Theorems 2 and 4:

$$\begin{aligned} \frac{1}{2}\text{-regret}(\text{Algorithm 2 applied on Algorithm 5}) &\leq O(n\sqrt{T \log n}). \\ \frac{1}{2}\text{-regret}(\text{Algorithm 4 applied on Algorithm 5}) &\leq O(n^{5/3} (\log n)^{1/3} T^{2/3}). \end{aligned}$$

This completes the proof. ■

## G.2. Proof of Corollary 1

*Proof.* The proof for the model from Asadpour et al. (2020) is a direct application of Theorem 5 by taking  $\lambda_i \triangleq \mathbb{P}_{u \sim \mathcal{G}}(\theta_u = i)$ , the probability that a consumer has patience level  $i$ , and  $f_i(S) \triangleq \mathbb{E}_{u \sim \mathcal{G}}[\kappa_u(S) | \theta_u = i]$ ,

the expected probability that a consumer with patience level  $i$  clicks on any of the top  $i$  products in  $S$ , as mentioned in Section 6.1. Thus, the sequential submodular function of interest is the expected probability that a consumer clicks on at least one product when offered an ordering  $\pi$  :

$$f(\pi) = \sum_{i=1}^n \lambda_i f_i(\{\pi(1), \dots, \pi(i)\}) = \sum_{i=1}^n \mathbb{P}_{u \sim \mathcal{G}}(\theta_u = i) \mathbb{E}_{e \sim \mathcal{G}}[\kappa_u(\pi) | \theta_u = i].$$

By invoking Theorem 5, we get the desired  $O(n\sqrt{T \log n})$   $\frac{1}{2}$ -regret in the full-information setting and  $O(n^{5/3}(\log n)^{1/3} T^{2/3})$   $\frac{1}{2}$ -regret in the bandit setting.

For the special consumer choice model in Ferreira et al. (2021), a consumer is characterized by two parameters: distribution of clicks for each item  $\mathbf{q}_u = (q_{u,1}, \dots, q_{u,n})$  and attention window size  $k_u$ . A consumer  $u$  examines the items in the top  $k_u$  positions, and an examined item  $i$  is clicked with probability  $q_{u,i}$  while unexamined items are never clicked. The events of clicking on two different items  $i$  or  $j$  in the event window are assumed to be independent. Notice that this is a special case of the choice model by Asadpour et al. (2020), where  $\theta_u = k_u$  and

$$\kappa_u(\{\pi(1), \dots, \pi(\theta_u)\}) = \kappa_u(\{\pi(1), \dots, \pi(k_u)\}) = 1 - \prod_{i=1}^{k_u} (1 - q_{u,i}).$$

The probability of click function  $\kappa_u$  is monotone since when  $X \subseteq Y \subseteq [n]$ , we have  $\prod_{i \in X} (1 - q_{u,i}) \geq \prod_{i \in Y} (1 - q_{u,i})$  (as  $0 \leq q_{u,i} \leq 1$  for all  $u$  and  $i$ ), which implies  $\kappa_u(X) \leq \kappa_u(Y)$ . It is also submodular, as for all  $X \subset Y \subseteq [n]$  and any item  $j \notin Y, j \in [n]$ , we have

$$\begin{aligned} & 1 - \prod_{i \in Y \setminus X} (1 - q_{u,i}) \geq 0 \\ \Leftrightarrow & (1 - (1 - q_j)) \left( \prod_{i \in X} (1 - q_{u,i}) \right) \left( 1 - \prod_{i \in Y \setminus X} (1 - q_{u,i}) \right) \geq 0 \\ \Leftrightarrow & \prod_{i \in X} (1 - q_{u,i}) - (1 - q_{u,j}) \prod_{i \in X} (1 - q_{u,i}) \geq \prod_{i \in Y} (1 - q_{u,i}) - (1 - q_{u,j}) \prod_{i \in Y} (1 - q_{u,i}) \\ \Leftrightarrow & \prod_{i \in X} (1 - q_{u,i}) - \prod_{i \in X \cup \{j\}} (1 - q_{u,i}) \geq \prod_{i \in Y} (1 - q_{u,i}) - \prod_{i \in Y \cup \{j\}} (1 - q_{u,i}) \\ \Leftrightarrow & \kappa_u(X \cup \{j\}) - \kappa_u(X) \geq \kappa_u(Y \cup \{j\}) - \kappa_u(Y). \end{aligned}$$

Since this choice model is a special case of the choice model in Corollary 1, we can invoke Corollary 1 to get the desired  $O(n\sqrt{T \log n})$   $\frac{1}{2}$ -regret in the full-information setting and  $O(n^{5/3}(\log n)^{1/3} T^{2/3})$   $\frac{1}{2}$ -regret in the bandit setting. ■

## Appendix H: Proofs and Remarks of Section 6.2 – Maximizing Multiple Reserves

In this appendix, first provide a discussion on major difference between Algorithm 6 and the algorithm in Roughgarden and Wang (2019). We then give the missing proofs of the results from Section 6.2. These results are restated for convenience.

## H.1. Discussion on Algorithm 6

The main difference between our algorithm and the algorithm in Roughgarden and Wang (2019) is the choice of revenue-from-reserves function  $q$ . Their revenue-from-reserves function  $q$  is different (coordinate-wise less) than ours. As it becomes more clear later in the proof, the need to design a new revenue-from-reserves function stems from our requirement to construct an explore sampling device for the online bandit learning algorithm.

## H.2. Proof of Theorem 6

*Proof.* We will show that our meta Algorithms 2 and 4 work by verifying the following conditions.

(i) Algorithm 6 is an extended  $(\frac{1}{2}, \frac{1}{2})$ -robust approximation algorithm. By Definition 6, we need to show that if each coordinate of our vector payoffs is bounded by some function  $h$ :

$$\forall j \in [m], \quad \left[ \sum_{t=1}^T \text{PAYOFF}^{(i)}(\tilde{\theta}_t^{(i)}, \mathbf{r}_t^{(i-1)}, \mathbf{v}_t) \right]_j \geq -h(T),$$

then we must have that our overall solution's error is bounded by:

$$\forall \mathbf{r}^* \in \mathcal{C}, \quad \sum_{t=1}^T \mathbb{E}[f(\mathbf{r}_t, \mathbf{v}_t)] \geq \frac{1}{2} \cdot \sum_{t=1}^T f(\mathbf{r}^*, \mathbf{v}_t) - \frac{1}{2}nh(T).$$

Recall from Section 6.2, that we defined the  $j^{\text{th}}$  coordinate of this vector payoff to be ( $j \in [m]$ ),  $\left[ \text{PAYOFF}^{(i)}(\theta^{(i)}, \mathbf{r}^{(i-1)}, \mathbf{v}) \right]_j \triangleq \mathbb{E}_{z' \sim \theta^{(i)}} [q^{(i)}(z') - q^{(i)}(\rho_j)]$ , and that  $\mathbf{r}_t^{(i)}$  is the reserve vector after subproblem  $i$ . Let's now define  $S_i$  to be the set of rounds where bidder  $i$  has the highest bid. We now carry out the standard offline analysis, but summed over all rounds  $t \in [T]$ .

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[f(\mathbf{r}_t, \mathbf{v}_t)] &\stackrel{(1)}{=} \frac{1}{2} \sum_{t=1}^T [\mathbf{v}_t]_{\hat{j}_t} + \frac{1}{2} \mathbb{E} \left[ \sum_{i=1}^n \sum_{t \in S_i} [\mathbf{r}_t]_i \mathbb{1} \left[ [\mathbf{r}_t]_i \in [[\mathbf{v}_t]_{\hat{j}_t}, [\mathbf{v}_t]_{\hat{j}_t^*}] \right] \right] \\ &\stackrel{(2)}{=} \frac{1}{2} \sum_{t=1}^T [\mathbf{v}_t]_{\hat{j}_t} + \frac{1}{2} \mathbb{E} \left[ \sum_{i=1}^n \sum_{t=1}^T q_t^{(i)}([\mathbf{r}_t]_i) \right] \\ &= \frac{1}{2} \sum_{t=1}^T [\mathbf{v}_t]_{\hat{j}_t} + \frac{1}{2} \sum_{i=1}^n \sum_{t \in S_i} q_t^{(i)}([\mathbf{r}^*]_i) + \frac{1}{2} \mathbb{E} \left[ \sum_{i=1}^n \sum_{t=1}^T q_t^{(i)}([\mathbf{r}_t]_i) \right] - \frac{1}{2} \sum_{i=1}^n \sum_{t \in S_i} q_t^{(i)}([\mathbf{r}^*]_i) \\ &\stackrel{(3)}{\geq} \frac{1}{2} \sum_{t=1}^T f(\mathbf{r}^*, \mathbf{v}_t) + \frac{1}{2} \sum_{i=1}^n \sum_{t=1}^T \mathbb{E} \left[ q_t^{(i)}([\mathbf{r}_t]_i) \right] - \frac{1}{2} \sum_{i=1}^n \sum_{t \in S_i} q_t^{(i)}([\mathbf{r}^*]_i) \\ &\stackrel{(4)}{=} \frac{1}{2} \sum_{t=1}^T f(\mathbf{r}^*, \mathbf{v}_t) + \frac{1}{2} \sum_{i=1}^n \sum_{t=1}^T \left[ \text{PAYOFF}^{(i)}(\tilde{\theta}_t^{(i)}, \mathbf{z}_t^{(i-1)}, f_t) \right]_{b_i} \\ &\stackrel{(5)}{\geq} \frac{1}{2} \sum_{t=1}^T f_t(\mathbf{z}^*) - \frac{1}{2}nh(T), \end{aligned}$$

where  $\hat{j}_t^*$  and  $\hat{j}_t$  are respectively the highest and second highest bidders in the valuation profile  $\mathbf{v}_t$ . We also defined a function  $q_t^{(i)}$  for each round, which is the same as the original  $q^{(i)}$  except with  $\mathbf{v}$  replaced by  $\mathbf{v}_t$ . Note that Inequality (1) holds because in each round  $t$ , with probability  $1/2$ , the algorithm returns  $\mathbf{z}_t = \mathbf{0}_n$ , which implies  $f(\mathbf{r}_t, \mathbf{v}_t) = [\mathbf{v}_t]_{\hat{j}_t}$  and with the same probability, it returns  $\mathbf{r}_t = \mathbf{r}_t^{(n)}$  which implies  $f(\mathbf{r}_t, \mathbf{v}_t)$  is at least equal to the reserve of the buyer with the highest bid when his reserve is less than his bid. Inequality (2) follows from the definition of  $q_t^{(i)}$ . Inequality (3) holds because under the optimal reserve price  $\mathbf{r}^*$ ,  $f(\mathbf{r}^*, \mathbf{v}_t)$

is less than or equal to the second highest bid  $[\mathbf{v}_t]_{j_t}$  when the bidder with the highest bid does not win or they win and their reserve is less than or equal to  $[\mathbf{v}_t]_{j_t}$ ; otherwise,  $f(\mathbf{r}^*, \mathbf{v}_t)$  is equal to  $q_t^{(i)}([\mathbf{r}_t]_{j_t^*}) \geq [\mathbf{v}_t]_{j_t}$ . Equality (4) follows from the definition of  $\text{PAYOFF}^{(i)}$ . In this inequality  $b_i$  is the index of element  $[\mathbf{r}^*]_i$ ; that is,  $[\mathbf{r}^*]_i = \rho_{b_i}$ . Recall that  $\tilde{\theta}_t^{(i)}$  is the (approximately-locally-optimal) distribution from which we are drawing  $[\mathbf{r}_t]_i$ . Finally, inequality (5) follows from the assumption. This inequality is the desired result.

(ii) *Algorithm 6 is bandit Blackwell reducible.* Per Definition 9, to show this statement, we will verify the following conditions:

- *Algorithm 6 is Blackwell reducible.* For every subproblem  $i \in [n]$ , consider an instance  $(\mathcal{X}, \mathcal{Y}, \mathbf{p}^{(i)})$  of Blackwell where  $\mathcal{X} = \Theta = \Delta(\mathcal{R})$  and  $\mathcal{Y} = [0, 1]^{d_{\text{param}}}$ , where  $d_{\text{param}} = |\mathcal{R}| = m$ . We can use the Blackwell adversary function (note that we identify adversary functions with valuation vector)  $\text{AdvB}^{(i)}(\mathbf{r}, \mathbf{v}) = [q^{(i)}(\rho_j)]_{j=1,2,\dots,m}$ .

The biaffine Blackwell payoff is  $\mathbf{p}^{(i)}(\boldsymbol{\theta}, \mathbf{y}) = \boldsymbol{\theta}^T \mathbf{y} \mathbf{1}_n - \mathbf{y}$  where  $\mathbf{1}_n$  is an  $n$ -dimensional all ones vector. Notice that the target set  $S$ , the non-negative orthant, is response-satisfiable because if player 1 plays  $\boldsymbol{\theta} = \mathbf{e}_{j^*}$  where  $j^* = \arg \max_{j \in [m]} [\mathbf{y}]_j$  then for every adversary's action  $\mathbf{y} \in \mathcal{Y}$ ,  $\mathbf{p}^{(i)}(\boldsymbol{\theta}, \mathbf{y}) \geq 0$ .

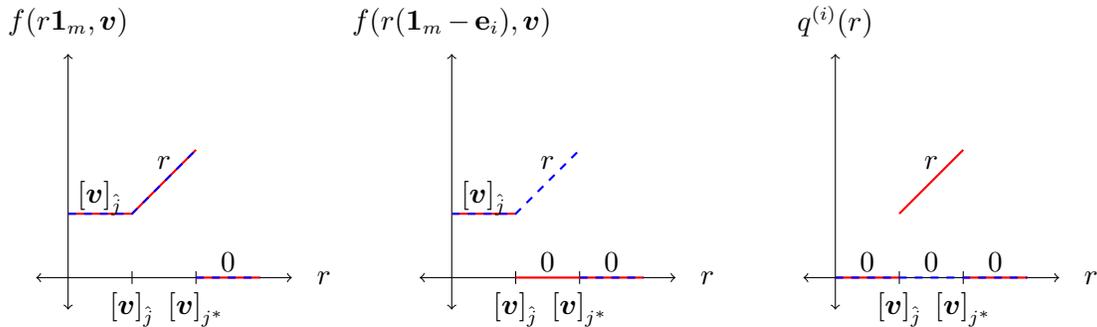
- *An unbiased estimator for the Blackwell payoff function  $\mathbf{p}^{(i)}$  can be constructed.* We will show that for every subproblem  $i \in [n]$ , there exists an explore sampling device  $U^{(i)}$  that returns  $(\mathbf{w}_{\text{exp}}^{(i)}, \mathbf{r}_{\text{exp}}^{(i)})$  such that (i) for all  $f \in \mathcal{F}, \boldsymbol{\theta} \in \Theta, \mathbf{r} \in \mathcal{D}, i \in [n]$ :  $\hat{\mathbf{p}}^{(i)}(\boldsymbol{\theta}, \text{AdvB}^{(i)}(\mathbf{r}, \mathbf{v})) = f(\mathbf{r}_{\text{exp}}^{(i)}, \mathbf{v}) \mathbf{w}_{\text{exp}}^{(i)}$ , where  $(\mathbf{w}_{\text{exp}}^{(i)}, \mathbf{r}_{\text{exp}}^{(i)}) \sim \text{ExpS}^{(i)}(\boldsymbol{\theta}, \mathbf{r})$ , and (ii)  $\hat{\mathbf{p}}$  is an unbiased estimator for the actual payoff, i.e.,  $\forall \boldsymbol{\theta} \in \Theta, \mathbf{y} \in \mathcal{Y} : \mathbb{E}[\hat{\mathbf{p}}^{(i)}(\boldsymbol{\theta}, \mathbf{y})] = \mathbf{p}^{(i)}(\boldsymbol{\theta}, \mathbf{y})$ . More specifically, we will construct a exploring distribution  $U^{(i)}$  such that if  $\mathbf{y} = \text{AdvB}(\mathbf{r}, f)$  for some  $f \in \mathcal{F}, \mathbf{r} \in \mathcal{D}$ , then  $\mathbb{E}[\hat{\mathbf{p}}(\boldsymbol{\theta}, \mathbf{y})] = \mathbb{E}[f(\mathbf{r}_{\text{exp}}) \mathbf{w}_{\text{exp}}] = \mathbf{p}(\boldsymbol{\theta}, \mathbf{y})$ , where the expectation is taken with respect to  $U^{(i)}$ . Notice that in Definition 9,  $U$  is not indexed by subproblems, but since the AdvB for this particular problem is subproblem specific, the distribution  $U$  should also depend on the subproblem. Because we would like to construct an unbiased estimator of the actual payoff  $\mathbf{p}^{(i)}$ , which is an affine function of  $\mathbf{y} = q^{(i)}$ , we focus on constructing an unbiased estimator for the function  $q^{(i)}$ . To do so, we make use of the following representation of  $q^{(i)}$ :

$$q^{(i)}(r) = f(r \mathbf{1}_n, \mathbf{v}) - f(r(\mathbf{1}_n - \mathbf{e}_i), \mathbf{v}). \quad (40)$$

To see what the above equation holds note that when bidder  $i$  does not have the highest bid in an auction, both  $q^{(i)}(r)$  and  $f(r \mathbf{1}_n, \mathbf{v}) - f(r(\mathbf{1}_n - \mathbf{e}_i), \mathbf{v})$ , which is the revenue gain of increasing bidder  $i$ 's reserve price from zero to  $r$ , are both zero. When bidder  $i$  has the highest bid in the auction,  $q^{(i)}(r) = r$  if  $r \in [[\mathbf{v}]_j, [\mathbf{v}]_{j^*}]$  and zero otherwise. Furthermore, the revenue from the reserve price  $r \mathbf{1}_n$ , i.e.,  $f(r \mathbf{1}_n, \mathbf{v})$ , is  $[\mathbf{v}]_j$  if  $r < [\mathbf{v}]_j$ ;  $r$  if  $r \in [[\mathbf{v}]_j, [\mathbf{v}]_{j^*}]$  and zero otherwise (the case  $r > [\mathbf{v}]_{j^*}$ ). The revenue from the reserve price  $r(\mathbf{1}_n - \mathbf{e}_i)$ , i.e.,  $f(r(\mathbf{1}_n - \mathbf{e}_i), \mathbf{v})$ , is  $[\mathbf{v}]_j$  if  $r < [\mathbf{v}]_j$  and zero otherwise. Thus,  $f(r \mathbf{1}_n, \mathbf{v}) - f(r(\mathbf{1}_n - \mathbf{e}_i), \mathbf{v})$  is  $r$  if  $r \in [[\mathbf{v}]_j, [\mathbf{v}]_{j^*}]$  and zero otherwise, which is exactly  $q^{(i)}(r)$ . This interesting relationship is depicted in Figure 2.

We now define the sampling distribution  $U^{(i)} : \mathcal{D} \times \Theta \rightarrow \Delta([0, 1]^m \times \mathcal{C})$ . For each  $j \in [m]$ , we pick:

$$\begin{aligned} (\mathbf{w}_{\text{exp}}^{(i)}, \mathbf{r}_{\text{exp}}^{(i)}) &= (2m(\boldsymbol{\theta}_j \mathbf{1}_n - \mathbf{e}_j), \rho_j \mathbf{1}_n), \quad \text{or} \\ (\mathbf{w}_{\text{exp}}^{(i)}, \mathbf{r}_{\text{exp}}^{(i)}) &= (-2m(\boldsymbol{\theta}_j \mathbf{1}_n - \mathbf{e}_j), \rho_j(\mathbf{1}_n - \mathbf{e}_i)), \end{aligned}$$



**Figure 2** The function  $q^{(i)}$  (right) and the two functions we combine to get it (left, center). The solid red line denotes the function value when  $i$  is the highest bidder, and the dashed blue line denotes the function value when  $i$  is not the highest bidder.

with probability  $\frac{1}{2m}$  each. Recall that  $\mathcal{D} = \mathcal{C} = \mathcal{R}^n$ , where  $\mathcal{R} = \{\rho_1, \rho_2, \dots, \rho_m\}$  is the set of possible reserve prices and  $\rho_j$  is the  $j$ -th largest reserve price in set  $\mathcal{R}$ . We then have

$$\begin{aligned}
\mathbb{E}[\hat{\mathbf{p}}^{(i)}(\boldsymbol{\theta}, \mathbf{y})] &= \mathbb{E}[\hat{\mathbf{p}}^{(i)}(\boldsymbol{\theta}, \text{AdvB}^{(i)}(r, \mathbf{v}))] = \mathbb{E}[f(\mathbf{r}_{\text{exp}}^{(i)}, \mathbf{v})\mathbf{w}_{\text{exp}}^{(i)}] \\
&= \boldsymbol{\theta}^T \begin{bmatrix} f(\rho_1 \mathbf{1}_n, \mathbf{v}) - f(\rho_1(\mathbf{1}_n - \mathbf{e}_i), \mathbf{v}) \\ \vdots \\ f(\rho_m \mathbf{1}_n, \mathbf{v}) - f(\rho_m(\mathbf{1}_n - \mathbf{e}_i), \mathbf{v}) \end{bmatrix} - \begin{bmatrix} f(\rho_1 \mathbf{1}_n, \mathbf{v}) - f(\rho_1(\mathbf{1}_n - \mathbf{e}_i), \mathbf{v}) \\ \vdots \\ f(\rho_m \mathbf{1}_n, \mathbf{v}) - f(\rho_m(\mathbf{1}_n - \mathbf{e}_i), \mathbf{v}) \end{bmatrix} \\
&= \boldsymbol{\theta}^T \begin{bmatrix} q^{(i)}(\rho_1) \\ \vdots \\ q^{(i)}(\rho_m) \end{bmatrix} - \begin{bmatrix} q^{(i)}(\rho_1) \\ \vdots \\ q^{(i)}(\rho_m) \end{bmatrix} = \boldsymbol{\theta}^T \mathbf{y} \mathbf{1}_n - \mathbf{y} = \mathbf{p}^{(i)}(\boldsymbol{\theta}, \mathbf{y}).
\end{aligned}$$

Wrapping up, Algorithm 6 is an extended  $(\frac{1}{2}, \frac{1}{2})$ -robust approximation algorithm with  $n$  subproblems and with a payoff diameter  $D(\mathbf{p})$  of  $O(1)$  and a payoff estimator diameter  $D(\hat{\mathbf{p}})$  of  $O(m)$ . It is also bandit Blackwell reducible. Therefore, from Theorems 2 and 4:

$$\begin{aligned}
\frac{1}{2}\text{-regret}(\text{Algorithm 2 applied on Algorithm 6}) &\leq O(nT^{1/2} \log^{1/2} m) \\
\frac{1}{2}\text{-regret}(\text{Algorithm 4 applied on Algorithm 6}) &\leq O(nm^{2/3} T^{2/3} \log^{1/3} m).
\end{aligned}$$

This completes the proof. ■

### H.3. Proof of Corollary 2

*Proof.* Let  $m \in \mathbb{Z}_+$  be a parameter we choose later to balance terms. We invoke Theorem 6 with the discretization  $\mathcal{R} = \{0, \frac{1}{m}, \frac{2}{m}, \dots, 1\}$ . Given any reserves  $\mathbf{r}^* \in [0, 1]^n$ , we can produce rounded reserves  $\tilde{\mathbf{r}}^*$  defined by rounding every reserve down to the nearest multiple of  $\frac{1}{m}$ :  $[\tilde{\mathbf{r}}]_i = \frac{1}{m} \lfloor m[r]_i \rfloor$ . Importantly, this never causes any bidder to fail to clear their reserve price (this is why we must round down and cannot round up). Hence this can only grow the set of bidders that clear their reserve and hence the maximum bid from this set can only increase. If a bidder that was already in this set proceeds to win the auction, then they are only competing with more bidders and their reserve price drops by at most  $1/m$ , so their payment can only drop by at most  $1/m$ . If a bidder not previously in this set proceeds to win the auction, then their

reserve price used to be higher than their valuation, but their valuation must be higher than the previous winner's valuation. They pay at least their reserve less  $\frac{1}{m}$ , so the revenue of the auction drops by at most  $1/m$  in this case as well. Hence the (summed) discretization error is  $T\frac{1}{m}$ , and we choose either  $m = \frac{1}{n}T^{1/2}$  (full-information) or  $m = n^{-3/5}T^{1/5}$  (bandit) to obtain:

$$\begin{aligned} O(nT^{1/2} \log^{1/2} m) + T\frac{1}{m} &= O(nT^{1/2} \log^{1/2} T) \\ O(nm^{2/3}T^{2/3} \log^{1/3} m) + T\frac{1}{m} &= O(n^{3/5}T^{4/5} \log^{1/3}(nT)) \end{aligned}$$

This completes the proof. ■

## Appendix I: Proofs and Remarks of Section E – Non-monotone Submodular Maximization

In this appendix, we first discuss the differences between Algorithm 7 and the bi-greedy algorithm by Niazadeh et al. (2018), and show that despite these differences Algorithm 7 obtains the same approximation factor as that of the bi-greedy algorithm. We then present the proof of Theorem 8.

### I.1. Discussion on Algorithm 7

Algorithm 7 is a modification of the bi-greedy algorithm by Niazadeh et al. (2018). But, as we show in this section, these modifications do not change the  $1/2$  approximation factor of the bi-greedy algorithm. We modify the bi-greedy algorithm to better satisfy the form of Algorithm 1, ease our construction of the sampling device and unbiased estimators in the bandit case, and provide a unified framework for submodular functions with a more general domain. The major differences and their corresponding reasons are as follows:

- To cover a more general discrete function domain, we optimize over points in the discrete set  $\mathcal{R}$  while their algorithm optimizes over  $[0, 1]^n$  implemented by casting an  $\epsilon$ -net.
- To help us construct the sampling device and unbiased estimators for the bandit case, in our local optimization step, we use  $\zeta^{(i)}(\hat{z}, z')$ , which is a linear combination of marginal functions  $\alpha^{(i)}$  and  $\beta^{(i)}$ , rather than  $\max\{\alpha^{(i)}(\hat{z}) - \alpha^{(i)}(z'), \beta^{(i)}(\hat{z}) - \beta^{(i)}(z')\}$ , in quantifying the value decrease of  $\hat{z}$ . Recall that in this step, we choose  $\theta^{(i)} \in \Delta(\mathcal{R})$  so that for all  $\hat{z} \in \mathcal{R}$ ,  $\mathbb{E}_{z' \sim \theta^{(i)}} [\frac{1}{2}\alpha^{(i)}(z') + \frac{1}{2}\beta^{(i)}(z') - \zeta^{(i)}(\hat{z}, z)] \geq 0$ .

Using the technique in Niazadeh et al. (2018), as we argue next, we can still find  $\theta^{(i)}$  that satisfies the condition in the local optimization step. Note that the bi-greedy analysis in Niazadeh et al. (2018) proves that satisfying this condition implies that Algorithm 7 is a  $\frac{1}{2}$ -approximation algorithm for the discretized submodular maximization problem.

**Satisfying the Local Optimization Step.** Here, we show how to choose  $\theta^{(i)} \in \Delta(\mathcal{R})$  that satisfies the condition in the local optimization step of Algorithm 7. To do so, First, we choose  $z_\ell \in \arg \max_{z \in \mathcal{R}} f(z, \bar{z}^{(i-1)})$  and  $z_u \in \arg \max_{z \in \mathcal{R}} f(z, \underline{z}^{(i-1)})$ . Then, we look at these two cases.

*Case (i):*  $z_u \leq z_\ell$ . We want to prove that deterministically returning  $z_\ell$  ( $\theta^{(i)}$  puts all its weight on  $z_\ell$ ) suffices. The key realization is that in this case,  $z_u$  and  $z_\ell$  maximize the functions  $f(\cdot, \bar{\mathbf{z}}^{(i-1)})$  and  $f(\cdot, \underline{\mathbf{z}}^{(i-1)})$  respectively:

$$f(z_\ell, \bar{\mathbf{z}}^{(i-1)}) \geq f(z_u, \bar{\mathbf{z}}^{(i-1)}), \quad f(z_u, \underline{\mathbf{z}}^{(i-1)}) \geq f(z_\ell, \underline{\mathbf{z}}^{(i-1)}). \quad (41)$$

We know by submodularity that two points are better than their coordinate-wise max and min:

$$f(z_u, \underline{\mathbf{z}}^{(i-1)}) + f(z_\ell, \bar{\mathbf{z}}^{(i-1)}) \leq f(z_\ell, \underline{\mathbf{z}}^{(i-1)}) + f(z_u, \bar{\mathbf{z}}^{(i-1)}). \quad (42)$$

Since adding up the first two inequalities in Equation (41) yields the third inequality in Equation (42), but with the direction reversed, we know all three must hold with equality. We conclude by noting that since  $z_\ell$  maximizes both functions, it also maximizes both  $\alpha^{(i)}$  and  $\beta^{(i)}$  at some nonnegative value and hence satisfies the desired condition in the local step optimization.

*Case (ii):*  $z_\ell < z_u$ . Suppose that the algorithm is able to find a  $\theta^{(i)}$  such that for any  $\hat{z} \in [z_\ell, z_u]$ , we have

$$\mathbb{E}_{z' \sim \theta^{(i)}} \left[ \frac{1}{2} \alpha^{(i)}(z') + \frac{1}{2} \beta^{(i)}(z') - \zeta^{(i)}(\hat{z}, z') \right] = \frac{1}{2} \alpha^{(i)}(z') + \frac{1}{2} \beta^{(i)}(z') - \zeta^{(i)}(\hat{z}, z') \geq 0. \quad (43)$$

We claim that this equation is still true for  $\hat{z}$  outside of the interval  $[z_\ell, z_u]$ .

Suppose that  $\hat{z} < z_\ell$ . By the choice of  $z_\ell$ , we know that  $\beta^{(i)}(z_\ell) \geq \beta^{(i)}(\hat{z})$ . By submodularity, we know that:

$$\begin{aligned} f(\hat{z}, \underline{\mathbf{z}}^{(i-1)}) + f(z_\ell, \bar{\mathbf{z}}^{(i-1)}) &\leq f(z_\ell, \underline{\mathbf{z}}^{(i-1)}) + f(\hat{z}, \bar{\mathbf{z}}^{(i-1)}) \\ \alpha^{(i)}(\hat{z}) + \beta^{(i)}(z_\ell) &\leq \alpha^{(i)}(z_\ell) + \beta^{(i)}(\hat{z}) \\ \beta^{(i)}(z_\ell) - \beta^{(i)}(\hat{z}) &\leq \alpha^{(i)}(z_\ell) - \alpha^{(i)}(\hat{z}). \end{aligned}$$

Since the LHS is nonnegative by the choice of  $z_\ell$ , so is the RHS. We have shown that  $z_\ell$  has strictly larger  $\alpha^{(i)}$  and  $\beta^{(i)}$  values (than  $\hat{z}$ ) and hence inequality in Equation (43) must be valid for  $\hat{z} < z_\ell$  as well. Analogous reasoning shows the same for the  $z_r < \hat{z}$  case. Notice that the method in Niazadeh et al. (2018) is able to compute a  $\theta^{(i)}$  that guarantees  $\mathbb{E}_{z' \sim \theta^{(i)}} \left[ \frac{1}{2} \alpha^{(i)}(z') + \frac{1}{2} \beta^{(i)}(z') - \zeta^{(i)}(\hat{z}, z') \right] \geq 0$  for any  $\hat{z} \in [z_\ell, z_u]$ , which means this is also true for any  $\hat{z} \in \mathcal{R}$ . Recall that the payoff function is

$$[\text{PAYOFF}(\theta, \mathbf{z}^{(i-1)}, f)]_j = \mathbb{E}_{z' \sim \theta} \left[ \frac{1}{2} \alpha^{(i)}(z') + \frac{1}{2} \beta^{(i)}(z') - \zeta^{(i)}(\rho_j, z') \right],$$

so such  $\theta^{(i)}$  also guarantees that  $\text{PAYOFF}(\theta^{(i)}, \mathbf{z}^{(i-1)}, f)$  is in the positive orthant.