

---

# A Bandit Learning Algorithm and Applications to Auction Design

---

Nguyễn Kim Thăng

IBISC, Univ. Evry, University Paris-Saclay, France  
kimthang.nguyen@univ-evry.fr

## Abstract

We consider online bandit learning in which at every time step, an algorithm has to make a decision and then observe only its reward. The goal is to design efficient (polynomial-time) algorithms that achieve a total reward approximately close to that of the best fixed decision in hindsight. In this paper, we introduce a new notion of  $(\lambda, \mu)$ -concave functions and present a bandit learning algorithm that achieves a performance guarantee which is characterized as a function of the concavity parameters  $\lambda$  and  $\mu$ . The algorithm is based on the mirror descent algorithm in which the update directions follow the gradient of the multilinear extensions of the reward functions. The regret bound induced by our algorithm is  $\tilde{O}(\sqrt{T})$  which is nearly optimal.

We apply our algorithm to auction design, specifically to welfare maximization, revenue maximization, and no-envy learning in auctions. In welfare maximization, we show that a version of fictitious play in smooth auctions guarantees a competitive regret bound which is determined by the smooth parameters. In revenue maximization, we consider the simultaneous second-price auctions with reserve prices in multi-parameter environments. We give a bandit algorithm which achieves the total revenue at least  $1/2$  times that of the best fixed reserve prices in hindsight. In no-envy learning, we study the bandit item selection problem where the player valuation is submodular and provide an efficient  $1/2$ -approximation no-envy algorithm.

## 1 Introduction

In Online Learning, the goal is to design algorithms which are robust in dynamically evolving environments by applying optimization methods that learn from experience and observations. Characterizing conditions, or in general discovering the hidden regularity, under which efficient online learning algorithms with performance guarantee exist is a major research agenda in online learning. In this paper, we consider this line of research and present a new regularity condition for the design of efficient online learning algorithms. Subsequently, we establish the applicability of our approach in auction design.

### 1.1 Definitions

**General online problem.** At each time step  $t = 1, 2, \dots$ , an algorithm chooses  $\mathbf{x}^t \in [0, 1]^n$ . After the algorithm has committed to its choice, an adversary selects a function  $f^t : [0, 1]^n \rightarrow [0, 1]$  that subsequently induces the reward of  $f^t(\mathbf{x}^t)$  for the algorithm. In the problem, we are interested in the bandit setting that at every time  $t$ , the algorithm observes only its reward  $f^t(\mathbf{x}^t)$ . The goal is to efficiently achieve the total gain approximately close to that obtained by the best decision in hindsight.

We consider the following notion of regret which measures the performance of algorithms.

**Definition 1** An algorithm is  $(r, R(T))$ -regret if for arbitrary total number of time steps  $T$  and for any sequence of reward functions  $f^1, \dots, f^T \in \mathcal{F}$ ,

$$\sum_{t=1}^T f^t(\mathbf{x}^t) \geq r \cdot \max_{\mathbf{x} \in [0,1]^n} \sum_{t=1}^T f^t(\mathbf{x}) - R(T).$$

We also say that the algorithm achieves a  $r$ -regret bound of  $R(T)$ .

In general, one seeks algorithms with  $(r, R(T))$ -regret such that  $r > 0$  is as large as possible (ideally, close to 1) and  $R(T)$  is sublinear as a function of  $T$ , i.e.,  $R(T) = o(T)$ . We also call  $r$  as the *approximation ratio* of the algorithm.

We introduce a regularity notion that generalizes the notion of concavity. The new notion, while simple, is crucial in our framework in order to design efficient online learning algorithms with performance guarantee.

**Definition 2** A function  $F$  is  $(\lambda, \mu)$ -concave if for every vectors  $\mathbf{x}$  and  $\mathbf{x}^*$ ,

$$\langle \nabla F(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle \geq \lambda F(\mathbf{x}^*) - \mu F(\mathbf{x}) \quad (1)$$

Note that a concave function is  $(1, 1)$ -concave. A non-trivial example is the  $(1, 2)$ -concavity of the multilinear relaxation of a monotone submodular function (Lemma 9).

## 1.2 Contribution

We aim to design a bandit algorithm for the general online problem with emphasis on auctions. Bandit algorithms have been widely studied in online convex optimization [23] but in the context of auction design, standard approaches have various limits. The main issues are: (1) the non-concavity of the (reward) functions, and (2) the intrinsic nature of the bandit setting (only the value  $f^t(\mathbf{x}^t)$  is observed). We overcome these issues by the approach which consists of lifting the search space and the reward functions to a higher dimension space and considering the multilinear extensions of the reward functions in that space. Concretely, we consider a sufficiently dense lattice  $\mathcal{L}$  in  $[0, 1]^n$  such that every point in  $[0, 1]^n$  can be approximated by a point in  $\mathcal{L}$ . Then, we lift all lattice points in  $\mathcal{L}$  to vertices of a hypercube in a high dimension space. Subsequently, we consider the multilinear extensions of reward functions  $f^t$  in that space. This procedure enables several useful properties, in particular the  $(\cdot, \cdot)$ -concavity, that hold for the multilinear extensions but not necessarily for the original reward functions. (For example, the multilinear extension of a monotone submodular function is always  $(1, 2)$ -concave but the submodular function is not.) The introduction of  $(\cdot, \cdot)$ -concavity and the use of multilinear extensions constitute the novel points in our approach compared to the previous ones. This allows us to bound the regret of our algorithm which is based on the classic mirror descent with respect to the gradients of the multilinear extensions.

**Informal Theorem 1** Let  $f^t : [0, 1]^n \rightarrow [0, 1]$  be the reward function at time  $1 \leq t \leq T$  and let  $F^t$  be the multilinear extension of the discretization of  $f^t$  on a lattice  $\mathcal{L}$ . Assume that  $f^t$ 's are  $L$ -Lipschitz and  $F^t$ 's are  $(\lambda, \mu)$ -concave. Then, there exists a bandit algorithm achieving

$$\sum_{t=1}^T \mathbb{E}[f^t(\mathbf{x}^t)] \geq \frac{\lambda}{\mu} \cdot \max_{\mathbf{x} \in [0,1]^n} \sum_{t=1}^T f^t(\mathbf{x}) - O(\max\{\lambda/\mu, 1\} Ln^{3/2}(\log T)^{3/2}(\log \log T)\sqrt{T}).$$

The formal statement corresponding to the above informal theorem is Theorem 2. By this theorem, determining the performance guarantee is reduced to computing the concave parameters. Moreover, the regret bound of  $\tilde{O}(\sqrt{T})$  is nearly optimal that has been proved in the context of online convex optimization (for concave functions, i.e.,  $(1, 1)$ -concave functions). The approach is convenient to derive bandit learning algorithms in the context of auction design as shown in the applications.

## Applications to Auction Design

In a general auction design setting, each player  $i$  has a *valuation* (or *type*)  $v_i$  and a set of actions  $\mathcal{A}_i$  for  $1 \leq i \leq n$ . Given an action profile  $\mathbf{a} = (a_1, \dots, a_n)$  consisting of actions chosen by players,

the auctioneer decides an allocation  $\mathbf{o}(\mathbf{a})$  and a payment  $p_i(\mathbf{o}(\mathbf{a}))$  for each player  $i$ . Note that for a fixed auction  $\mathbf{o}$ , the outcome  $\mathbf{o}(\mathbf{a})$  of the game is completely determined by the action profile  $\mathbf{a}$ . Then, the *utility* of player  $i$  with valuation  $v_i$ , following the quasi-linear utility model, is defined as  $u_i(\mathbf{o}(\mathbf{a}); v_i) = v_i(\mathbf{o}(\mathbf{a})) - p_i(\mathbf{o}(\mathbf{a}))$ . The *social welfare* of an auction is defined as the total utility of all participants (the players and the auctioneer):  $\text{SW}(\mathbf{o}(\mathbf{a}); \mathbf{v}) = \sum_{i=1}^n u_i(\mathbf{o}(\mathbf{a}); v_i) + \sum_{i=1}^n p_i(\mathbf{a})$ . The total revenue of the auction is  $\text{REV}(\mathbf{o}(\mathbf{a}), \mathbf{v}) = \sum_{i=1}^n p_i(\mathbf{o}(\mathbf{a}))$ . When  $\mathbf{o}$  is clear in the context, we simply write  $v_i(\mathbf{a}), u_i(\mathbf{a}; v_i), p_i(\mathbf{a}), \text{SW}(\mathbf{a}; \mathbf{v}), \text{REV}(\mathbf{a}, \mathbf{v})$  instead of  $v_i(\mathbf{o}(\mathbf{a})), u_i(\mathbf{o}(\mathbf{a}); v_i), p_i(\mathbf{o}(\mathbf{a})), \text{SW}(\mathbf{o}(\mathbf{a}); \mathbf{v}), \text{REV}(\mathbf{o}(\mathbf{a}), \mathbf{v})$ , respectively. In the paper, we consider two standard objectives: welfare maximization and revenue maximization. Note that in revenue maximization, we call players as bidders.

### 1.2.1 Fictitious Play in Smooth Auctions

We consider adaptive dynamics in auctions. In the setting, there is an underlying auction  $\mathbf{o}$  and there are  $n$  players, each player  $i$  has a set of actions  $\mathcal{A}_i$  and a valuation function  $v_i$  taking values in  $[0, 1]$  (by normalization). In every time step  $1 \leq t \leq T$ , each player  $i$  selects a strategy which is a distribution in  $\Delta(\mathcal{A}_i)$  according to some given adaptive dynamic. After all players have committed their strategies, which result in a strategy profile  $\boldsymbol{\sigma}^t \in \Delta(\mathcal{A})$ , the auction induces a social welfare  $\text{SW}(\mathbf{o}, \boldsymbol{\sigma}^t) := \mathbb{E}_{\mathbf{a} \sim \boldsymbol{\sigma}^t} [\text{SW}(\mathbf{o}(\mathbf{a}); \mathbf{v})]$ . In this setting, we study the total welfare achieved by the given adaptive dynamic comparing to the optimal welfare. This problem can be cast in the online optimization framework in which at time step  $t$ , the player strategy profile corresponds to the decision of the algorithm and subsequently, the gain of the algorithm is the social welfare induced by the auction w.r.t the strategy profile.

Smooth auctions is an important class of auctions in welfare maximization. The smoothness notion has been introduced [40, 36] in order to characterize the efficiency of (Bayes-Nash) equilibria of auctions. It has been shown that several auctions in widely studied settings are smooth; and many proof techniques analyzing equilibrium efficiency can be reduced to the smooth argument.

**Definition 3 ([40, 36])** For parameters  $\lambda, \mu \geq 0$ , an auction is  $(\lambda, \mu)$ -smooth if for every valuation profile  $\mathbf{v} = (v_1, \dots, v_n)$ , there exist action distributions  $\bar{D}_1(\mathbf{v}), \dots, \bar{D}_n(\mathbf{v})$  over  $\mathcal{A}_1, \dots, \mathcal{A}_n$  such that, for every action profile  $\mathbf{a}$ ,

$$\sum_{i=1}^n \mathbb{E}_{\bar{\mathbf{a}}_i \sim \bar{D}_i(\mathbf{v})} [u_i(\bar{\mathbf{a}}_i, \mathbf{a}_{-i}; v_i)] \geq \lambda \cdot \text{SW}(\bar{\mathbf{a}}; \mathbf{v}) - \mu \cdot \text{SW}(\mathbf{a}; \mathbf{v})$$

where  $\mathbf{a}_{-i}$  is the action profile similar to  $\mathbf{a}$  without player  $i$ .

It has been proved that if an auction is  $(\lambda, \mu)$ -smooth then every Bayes-Nash equilibrium of the auction has expected welfare at least  $\lambda/\mu$  fraction of the optimal auction [36, 40]. The performance guarantee holds even for vanishing regret sequences. A sequence of actions profiles  $\mathbf{a}^1, \mathbf{a}^2, \dots$ , is an *individually-vanishing-regret sequence* if for every player  $i$  and action  $a'_i$ ,

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T [u_i(a'_i, \mathbf{a}_{-i}^t; v_i) - u_i(\mathbf{a}^t; v_i)] \leq 0. \quad (2)$$

However, several interesting dynamics are not guaranteed to have the individually-vanishing-regret property. In a recent survey, Roughgarden et al. [38] have raised a question whether adaptive dynamics without the individually-vanishing-regret condition can achieve approximate optimal welfare. Among others, *fictitious play* [8] is an interesting, widely-studied dynamic which attracts a significant attention in the community.

In the paper, we consider a version of fictitious play in smooth auctions, namely Perturbed Discrete Time Fictitious Play (PDTFP). (The formal definition is given in Section 4.) In general, it is not known whether this dynamic has individually-vanishing-regret. (In particular, in the PDTFP dynamic players are valuation-oriented whereas the condition (2) concerns player utilities.) Despite that fact, using our framework, we prove that given an offline  $(\lambda, \mu)$ -smooth auction, PDTFP dynamic achieves a  $\lambda/(1 + \mu)$  fraction of the optimal welfare. The corresponding formal statement is Theorem 3.

**Informal Theorem 2** If the underlying auction  $\mathbf{o}$  is a  $(\lambda, \mu)$ -smooth then the PDTFP dynamic achieves  $\left(\frac{\lambda}{1+\mu}, R(T)\right)$ -regret where  $R(T) = O\left(\frac{\sqrt{T}}{1+\mu}\right)$ .

### 1.2.2 Revenue maximization in Multi-Dimensional Environments

We consider online simultaneous second-price auctions with reserve prices in *multi-dimensional* environments. In the setting, there are  $n$  bidders and  $m$  items to be sold to these bidders. At every time step  $t = 1, 2, \dots, T$ , the auctioneer selects reserve prices  $r_i^t = (r_{i1}^t, \dots, r_{im}^t)$  for each bidder  $i$  where  $r_{ij}^t$  is the reserve price of item  $j$  for bidder  $i$ . Each bidder  $i$  for  $1 \leq i \leq n$  has a (private) valuation  $v_i^t : 2^{[m]} \rightarrow \mathbb{R}^+$  over subsets of items. After the reserve prices have been chosen, every bidder  $i$  picks a bid vector  $b_i^t$  where  $b_{ij}^t$  is the bid of bidder  $i$  on item  $j$  for  $1 \leq j \leq m$ . Then the auction for each item  $1 \leq j \leq m$  works as follows: (1) remove all bidders  $i$  with  $b_{ij}^t < r_{ij}^t$ ; (2) run the second price auction on the remaining bidders to determine the winner of item  $j$  — the bidder with highest non-removed bid on item  $j$ ; and (3) charge the winner of item  $j$  the price which is the maximum of  $r_{ij}^t$  and the second highest bid among non-removed bids  $b_{ij}^t$ . The objective of the auctioneer is to achieve the total revenue approximately close to that achieved by the best fixed reserve-price auction. Note that in the bandit setting, the auction is given as a blackbox and at every time step, the auctioneer observes only the total revenue (total price) without knowing neither the bids of bidders nor the winner/price of each item. The setting enhances, among others, the privacy of bidders.

The second-price auctions with reserve prices in *single-parameter* environments have been considered by Roughgarden and Wang [37] in full-information online learning. Using the Follow-the-Perturbed-Leader strategy, they gave a polynomial-time online algorithm that achieves half the revenue of the best fixed reserve-price auction minus a term  $O(\sqrt{T} \log T)$  (so their algorithm is  $(1/2, O(\sqrt{T} \log T))$ -regret in our terminology). The problem we consider cannot be reduced to applying their algorithm on  $m$  separate items since (1) bids on different items might be highly correlated (due to bidders' valuations); and (2) in the bandit setting for multiple items, the auctioneer know only the total revenue (not the revenue from each item). Using our framework, we prove the following result.

**Informal Theorem 3** *There exist a bandit algorithm that achieves the regret bound of  $(1/2, O(mn^{3/2}(\log T)^{3/2}(\log \log T)\sqrt{T}))$  for online simultaneous second-price auctions with reserve prices in multi-parameter environments.*

### 1.2.3 Bandit No-Envy Learning in Auctions

The concept of *no-envy learning* in auctions has been introduced by Daskalakis and Syrgkanis [14] in order to maintain approximate welfare optimality while guaranteeing computational tractability. The concept is inspired by the notion of *Walrasian equilibrium*. Intuitively, an allocation of items to buyers together with a price on each item forms a Walrasian equilibrium if no buyer envies other allocation given the current prices. In the paper, we consider no-envy bandit learning algorithms for the following *online item selection* problem.

In the problem, there are  $m$  items and a player with monotone valuation  $v : 2^{[m]} \rightarrow \mathbb{R}^+$ . At every time step  $1 \leq t \leq T$ , the player chooses a subset of items  $S^t \subset [m]$  and the adversary picks adaptively (probably depending on the history up to time  $t - 1$  but not on the current set  $S^t$ ) a threshold vector  $\mathbf{p}^t$ . The player observed the total price  $\sum_{j \in S^t} p_j^t$  and gets the reward of  $v(S^t) - \sum_{j \in S^t} p_j^t$ . A learning algorithm for the online item selection problem is a *r-approximate no-envy* if for any adaptively chosen sequence of threshold vectors  $\mathbf{p}^t$  for  $1 \leq t \leq T$ , the sets  $S^t$  for  $1 \leq t \leq T$  chosen by the algorithm satisfy

$$\mathbb{E} \left[ \sum_{t=1}^T \left( v(S^t) - \sum_{j \in S^t} p_j^t \right) \right] \geq \max_{S \subset [m]} \sum_{t=1}^T \left( r \cdot v(S) - \sum_{j \in S} p_j^t \right) - R(T)$$

where the regret  $R(T) = o(T)$ .

Daskalakis and Syrgkanis [14] considered the problem in the full-information setting (i.e., at every time step  $t$ , the player observes the whole vector  $\mathbf{p}^t$ ) where the valuation  $v$  is a coverage function<sup>1</sup> and gave an  $(1 - 1/e)$ -approximate no-envy algorithm with regret bound  $O(\sqrt{T})$ . The algorithm is designed via the convex rounding scheme [16], a technique which has been used in approximation algorithms and in truthful mechanism design.

<sup>1</sup>A coverage function  $v : 2^{[m]} \rightarrow \mathbb{R}^+$  has the form  $v(S) = |\cup_{j \in S} A_j|$  where  $A_1, \dots, A_m$  are subsets of  $[m]$ .

In this paper, we consider *submodular* valuations, a more general and widely-studied class of valuations. A valuation  $v : 2^{[m]} \rightarrow \mathbb{R}^+$  is *submodular* if for any sets  $S \subset T \subset [m]$ , and for every item  $j$ , it holds that  $v(S \cup j) - v(S) \geq v(T \cup j) - v(T)$ . Using our framework, we prove the following result.

**Informal Theorem 4** *There exist an  $(1/2, O(m^{3/2}(\log T)^{3/2}(\log \log T)\sqrt{T}))$ -regret no-envy learning algorithm for the bandit item selection problem where the player valuation is submodular.*

### 1.3 Related Work

There is large literature on online learning and auction design. In this section, we summarize and discuss only works directly related to ours. The interested reader can refer to [39, 23] for online learning and to [38] (and references therein) for auction design.

**Online/Bandit Learning.** Online Learning, or Online Convex Optimization, is an active research domain. The first no-regret algorithm was given by Hannan [21]. Subsequently, Littlestone and Warmuth [30] and Freund and Schapire [18] gave improved algorithms with regret  $\sqrt{\log(|\mathcal{A}|)}o(T)$  where  $|\mathcal{A}|$  is the size of the action space. Kalai and Vempala [27] presented the first efficient online algorithm, called *Follow-the-Perturbed-Leader* (FTPL), for linear objective functions. The strategy consists of adding perturbation to the cumulative gain (payoff) of each action and then selecting the action with the highest perturbed gain. This strategy has been generalized and successfully applied to several settings [24, 41, 14, 15]. Specifically, FTPL and its generalized versions have been used to design efficient online no-regret algorithms with oracles beyond the linear setting: to submodular [24] and non-convex [2] settings.

In bandit learning, many interesting results and powerful optimization/algorithmic methods have been proved and introduced, including interior point methods [1], random walk [33], continuous multiplicative updates [13], random perturbation [3], iterative methods [17]. In bandit linear optimization, the near-optimal regret bound of  $\tilde{O}(n\sqrt{T})$  has been established due to a long line of works [1, 13, 10]. Beyond the linear functions, several results have been known. Kleinberg [29], Flaxman et al. [17] provided  $\tilde{O}(\text{poly}(n)T^{3/4})$ -regret algorithm for general convex functions. Subsequently, Hazan and Li [25] presented an (exponential-time) algorithm which achieves  $\tilde{O}(\exp(n)\sqrt{T})$ -regret. Recently, Bubeck et al. [11] gave the first polynomial-time algorithm with regret  $\tilde{O}(n^{9.5}\sqrt{T})$ .

**Smooth Auctions and Fictitious Play.** The smoothness framework was introduced in order to prove approximation guarantees for equilibria in complete-information [35] and incomplete-information [40, 36] games. Smooth auctions (Definition 3) is a large class of auctions where the price of anarchy can be systematically characterized by the smooth arguments. Many interesting auctions have been shown to be smooth; and the smooth argument is a central proof technique to analyze the price of anarchy. We refer the reader to a recent survey [38] for more details. The smoothness framework extends to adaptive dynamics with vanishing regret. However, several important dynamics are not guaranteed to have the vanishing regret property, for example the class of fictitious play [8] and other classes of dynamics in [20]. A research agenda, as raised in [38], is to characterize the performance of such dynamics. Some recent works (e.g., [31]) have been considered in this direction.

**Revenue Maximization.** Optimal truthful auctions in single-parameter environments are completely characterized by Myerson [32]. Recently, a major line of research in data-driven mechanism design focus on competitive auctions without the full knowledge on the valuation distribution and even in non-stochastic settings. The study of second-price auctions with reserve prices in single-parameter environments and its variants have been considered in [28, 7, 12]. Recently, Roughgarden and Wang [37] gave a polynomial-time online algorithm that achieves  $(1/2, O(\sqrt{T}))$ -regret. Subsequently, Dudik et al. [15] showed that the same regret bound can be obtained using their framework. Both are in the online full-information setting.

**No-envy Learning in Auctions.** The notion of *no-envy learning* in auctions has been introduced by Daskalakis and Syrgkanis [14]. They proposed the concept of no-envy learning in order to maintain both the welfare optimality and computational tractability. Among others, Daskalakis and Syrgkanis

[14] considered the online item selection problem with coverage valuation and gave an efficient  $(1 - 1/e)$ -approximate no-envy algorithm with regret bound of  $O(\sqrt{T})$ .

## 1.4 Organization

We begin by giving some preliminary definitions in Section 2. In Section 3, we present our framework and the main algorithm. Subsequently, we show the applications about: (1) Perturbed Fictitious Play in Smooth Auctions in Section 4; (2) Online Simultaneous Second-Price Auctions with Reserve Prices in Section 5; and (3) Bandit No-Envy Learning in Auctions in Section 6.

## 2 Preliminaries

Given a norm  $\|\cdot\|$ , the dual norm is defined as  $\|\mathbf{y}\|_* := \max_{\mathbf{x}: \|\mathbf{x}\|=1} \langle \mathbf{x}, \mathbf{y} \rangle$ . A function  $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\alpha_\Phi$ -strongly convex w.r.t  $\|\cdot\|$  if

$$\Phi(\mathbf{x}') \geq \Phi(\mathbf{x}) + \langle \nabla \Phi(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{\alpha_\Phi}{2} \|\mathbf{x}' - \mathbf{x}\|^2$$

Given a strongly convex function  $\Phi$ , a map defined as  $\mathbf{x} \mapsto \nabla \Phi(\mathbf{x})$  is bijective. Denote  $\nabla \Phi^*$  the inverse map of  $\nabla \Phi$ . In fact, this inverse map is given by the gradient of the Fenchel dual for  $\Phi$ . We refer reader to [5] and [6, Chapter 7] for more details. Given a strictly convex function  $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}$ , the *Bregman divergence* is defined as

$$D_\Phi(\mathbf{x} \|\mathbf{x}') := \Phi(\mathbf{x}) - \Phi(\mathbf{x}') - \langle \nabla \Phi(\mathbf{x}'), \mathbf{x} - \mathbf{x}' \rangle$$

The following lemma generalizes the Pythagoras theorem (proof can be found in [4] for example).

**Lemma 1 (Generalized Pythagorean Property)** *Given a convex body  $\mathcal{K} \subset \mathbb{R}^n$ . Let  $\mathbf{x} \in \mathcal{K}$  and  $\mathbf{y}' \in \mathbb{R}^n$ . Let  $\mathbf{y}$  be the projection of  $\mathbf{y}'$  on  $\mathcal{K}$ , defined as  $\mathbf{y} = \arg \min_{\bar{\mathbf{y}} \in \mathcal{K}} D_\Phi(\bar{\mathbf{y}} \|\mathbf{y}')$ . Then  $D_\Phi(\mathbf{x} \|\mathbf{y}) \leq D_\Phi(\mathbf{x} \|\mathbf{y}')$ .*

Let  $\mathcal{K} \subset \mathbb{R}^n$  be a convex set with non-empty interior  $\text{int}(\mathcal{K})$ . A function  $\Phi : \mathcal{K} \rightarrow \mathbb{R}$  is  $\nu$ -self-concordant if

1.  $\Phi$  is three times continuously differentiable, convex and  $\Phi$  approaches infinity along any sequence approaching the boundary of  $\mathcal{K}$
2. For every  $\mathbf{a} \in \mathcal{R}^n$  and  $\mathbf{x} \in \text{int}(\mathcal{K})$ , it holds that

$$\begin{aligned} |\nabla^3 \Phi(\mathbf{x})[\mathbf{a}, \mathbf{a}, \mathbf{a}]| &\leq 2(|\nabla^2 \Phi(\mathbf{x})[\mathbf{a}, \mathbf{a}]|)^{3/2} \\ |\nabla \Phi(\mathbf{x})[\mathbf{a}]| &\leq \nu^{1/2} (|\nabla^2 \Phi(\mathbf{x})[\mathbf{a}, \mathbf{a}]|)^{1/2} \end{aligned}$$

$$\text{where } \nabla^3 \Phi(\mathbf{x})[\mathbf{a}, \mathbf{a}, \mathbf{a}] := \frac{\partial^3}{\partial t_1 \partial t_2 \partial t_3} \Phi(\mathbf{x} + t_1 \mathbf{a} + t_2 \mathbf{a} + t_3 \mathbf{a}) \Big|_{t_1=t_2=t_3=0}$$

One can define a local norm based on the Hessian of a self-concordant function. Formally, given a  $\nu$ -self-concordant function  $\Phi$  and a point  $\mathbf{x} \in \text{int}(\mathcal{K})$ , the local norm  $\|\cdot\|_{\mathbf{x}}$  and its dual norm  $\|\cdot\|_{*,\mathbf{x}}$  are defined as

$$\|\mathbf{a}\|_{\mathbf{x}} = (\mathbf{a}^\top \nabla^2 \Phi(\mathbf{x}) \mathbf{a})^{1/2}, \quad \|\mathbf{a}\|_{*,\mathbf{x}} = (\mathbf{a}^\top (\nabla^2 \Phi(\mathbf{x}))^{-1} \mathbf{a})^{1/2}$$

The following lemma states a useful property of self-concordant functions.

**Lemma 2 ([23], Lemma 6.10)** *Let  $\Phi$  be a  $\nu$ -self-concordant function over a convex set  $\mathcal{K}$ . Then, for all  $\mathbf{x}, \mathbf{y} \in \text{int}(\mathcal{K})$ ,*

$$\Phi(\mathbf{x}) - \Phi(\mathbf{y}) \leq \nu \log \frac{1}{1 - \pi_{\mathbf{x}}(\mathbf{y})}$$

where  $\pi$  is the Minkowski function over  $\mathcal{K}$  defined as  $\pi_{\mathbf{x}}(\mathbf{y}) := \inf\{t \geq 0 : \mathbf{x} + t^{-1}(\mathbf{y} - \mathbf{x}) \in \mathcal{K}\}$ .

### 3 Framework of Online Learning

We present a general efficient online algorithm and characterize its regret bound based on its concavity parameters. In Section 3.1, we prove the guarantee of the online mirror descent algorithm assuming access to unbiased estimates of the gradients of the functions. In Section 3.2, we derive an algorithm in the bandit setting.

#### 3.1 Regret of $(\lambda, \mu)$ -Concave Functions

**Mirror descent.** Given a convex set  $\mathcal{K}$ . Let  $\Phi$  be a  $\alpha_\Phi$ -strongly convex function w.r.t  $\|\cdot\|$ . Initially, let  $\mathbf{x}^1$  is an arbitrary point in  $\mathcal{K}$ . At time step  $t$ , play  $\mathbf{x}^t$  and receive the reward of  $F^t(\mathbf{x}^t)$ . Let  $\mathbf{g}^t$  be an unbiased estimate of  $-\nabla F^t(\mathbf{x}^t)$  and denote  $\boldsymbol{\theta}^t = \nabla\Phi(\mathbf{x}^t)$ . The algorithm selects the decision  $\mathbf{x}^{t+1}$  as follows.

$$\begin{aligned}\zeta^{t+1} &= \boldsymbol{\theta}^t - \eta \cdot \mathbf{g}^t \\ \mathbf{y}^{t+1} &= \nabla\Phi^*(\zeta^{t+1}) \\ \mathbf{x}^{t+1} &= \arg \min_{\mathbf{x} \in \mathcal{K}} D_\Phi(\mathbf{x} \|\mathbf{y}^{t+1})\end{aligned}$$

where  $\eta$  is a step size. An equivalent description is

$$\mathbf{x}^{t+1} = \arg \max_{\mathbf{x} \in \mathcal{K}} \{\langle \eta \mathbf{g}^t, \mathbf{x} - \mathbf{x}^t \rangle - D_\Phi(\mathbf{x} \|\mathbf{x}^t)\}. \quad (3)$$

**Theorem 1** Assume that  $F^t$  is  $(\lambda, \mu)$ -concave for every  $1 \leq t \leq T$  and  $\mathbf{x}^*$  is the best solution in hindsight, i.e.,  $\mathbf{x}^* \in \arg \max_{\mathbf{x}} \sum_{t=1}^T F^t(\mathbf{x})$ . Then the mirror descent algorithm achieves  $(\frac{\lambda}{\mu}, R(T))$ -regret in expectation where

$$R(T) = \frac{1}{\mu \cdot \eta} D_\Phi(\mathbf{x}^* \|\mathbf{x}^1) + \frac{\eta}{\mu \cdot 2\alpha_\Phi} \sum_{t=1}^T \|\mathbf{g}^t\|_*^2$$

If  $\|\mathbf{g}^t\|_* \leq L_g$  for  $1 \leq t \leq T$  (i.e.,  $F^t$  is  $L_g$ -Lipschitz w.r.t  $\|\cdot\|$ ) and  $D_\Phi(\mathbf{x}^* \|\mathbf{x}^1)$  is bounded by  $G^2$  then by choosing  $\eta = \frac{G}{L_g} \sqrt{\frac{2\alpha_\Phi}{T}}$ , we have  $R(T) \leq \frac{GL_g}{\mu} \sqrt{2\alpha_\Phi T}$ .

*Proof* In the analysis, we follow the potential argument of Bansal and Gupta [4] and derive a bound based on the concavity parameters. Define the potential as  $\Psi^t = \frac{1}{\eta} D_\Phi(\mathbf{x}^* \|\mathbf{x}^t)$ . First, we observe that

$$\begin{aligned}\eta(\Psi^{t+1} - \Psi^t) &= D_\Phi(\mathbf{x}^* \|\mathbf{x}^{t+1}) - D_\Phi(\mathbf{x}^* \|\mathbf{x}^t) \\ &\leq D_\Phi(\mathbf{x}^* \|\mathbf{y}^{t+1}) - D_\Phi(\mathbf{x}^* \|\mathbf{x}^t) \\ &= \Phi(\mathbf{x}^*) - \Phi(\mathbf{y}^{t+1}) - \underbrace{\langle \nabla\Phi(\mathbf{y}^{t+1}), \mathbf{x}^* - \mathbf{y}^{t+1} \rangle}_{\zeta^{t+1}} - \Phi(\mathbf{x}^*) + \Phi(\mathbf{x}^t) + \underbrace{\langle \nabla\Phi(\mathbf{x}^t), \mathbf{x}^* - \mathbf{x}^t \rangle}_{\boldsymbol{\theta}^t} \\ &= \Phi(\mathbf{x}^t) - \Phi(\mathbf{y}^{t+1}) - \langle \zeta^{t+1}, \mathbf{x}^t - \mathbf{y}^{t+1} \rangle - \langle \zeta^{t+1} - \boldsymbol{\theta}^t, \mathbf{x}^* - \mathbf{x}^t \rangle \\ &= \Phi(\mathbf{x}^t) - \Phi(\mathbf{y}^{t+1}) - \langle \boldsymbol{\theta}^t, \mathbf{x}^t - \mathbf{y}^{t+1} \rangle + \langle \eta \mathbf{g}^t, \mathbf{x}^t - \mathbf{y}^{t+1} \rangle + \langle \eta \mathbf{g}^t, \mathbf{x}^* - \mathbf{x}^t \rangle \\ &\leq -\frac{\alpha_\Phi}{2} \|\mathbf{y}^{t+1} - \mathbf{x}^t\|^2 + \eta \langle \mathbf{g}^t, \mathbf{x}^t - \mathbf{y}^{t+1} \rangle + \eta \langle \mathbf{g}^t, \mathbf{x}^* - \mathbf{x}^t \rangle \\ &= -\frac{\alpha_\Phi}{2} \|\mathbf{y}^{t+1} - \mathbf{x}^t\|^2 + \frac{1}{\alpha_\Phi} \langle \eta \mathbf{g}^t, \alpha_\Phi (\mathbf{x}^t - \mathbf{y}^{t+1}) \rangle + \eta \langle \mathbf{g}^t, \mathbf{x}^* - \mathbf{x}^t \rangle \\ &\leq \frac{\eta^2}{2\alpha_\Phi} \|\mathbf{g}^t\|_*^2 + \eta \langle \mathbf{g}^t, \mathbf{x}^* - \mathbf{x}^t \rangle\end{aligned}$$

where the first inequality is due to the generalized Pythagorean property (Lemma 1); the fourth equality follows the update rule  $\zeta^{t+1} = \boldsymbol{\theta}^t - \eta \cdot \mathbf{g}^t$ ; the second inequality holds since  $\Phi$  is  $\alpha_\Phi$ -strongly convex; and in the last inequality, we use Cauchy-Schwarz inequality  $\langle \mathbf{a}, \mathbf{b} \rangle \leq \|\mathbf{b}\| \|\mathbf{a}\|_* \leq \|\mathbf{b}\|^2/2 + \|\mathbf{a}\|_*^2/2$ .

By the observation and the fact that  $\mathbf{g}^t$  is an unbiased estimate of  $-\nabla F^t(\mathbf{x}^t)$ ,

$$\mathbb{E}[(\Psi^{t+1} - \Psi^t)] \leq \frac{\eta}{2\alpha_\Phi} \mathbb{E}[\|\mathbf{g}^t\|_*^2] - \langle \nabla F^t(\mathbf{x}^t), \mathbf{x}^* - \mathbf{x}^t \rangle. \quad (4)$$

Using the bound of the potential change due to Inequality (4) and linearity of expectation, we get

$$\begin{aligned}
\mathbb{E} \left[ \sum_{t=1}^T (\lambda F^t(\mathbf{x}^*) - \mu F^t(\mathbf{x}^t)) \right] &\leq \Psi^1 + \sum_{t=1}^T \mathbb{E} \left[ \lambda F^t(\mathbf{x}^*) - \mu F^t(\mathbf{x}^t) + \Psi^{t+1} - \Psi^t \right] \\
&\leq \Psi^1 + \sum_{t=1}^T \mathbb{E} \left[ \underbrace{\lambda F^t(\mathbf{x}^*) - \mu F^t(\mathbf{x}^t) - \langle \nabla F^t(\mathbf{x}^t), \mathbf{x}^* - \mathbf{x}^t \rangle}_{\leq 0 \text{ since } F^t \text{ is } (\lambda, \mu)\text{-concave}} + \frac{\eta}{2\alpha_\Phi} \|\mathbf{g}^t\|_*^2 \right] \\
&\leq \frac{1}{\eta} D_\Phi(\mathbf{x}^* \|\mathbf{x}^1) + \frac{\eta}{2\alpha_\Phi} \sum_{t=1}^T \mathbb{E} [\|\mathbf{g}^t\|_*^2]. \tag{5}
\end{aligned}$$

If the norms  $\|\mathbf{g}^t\|_*$  are bounded by  $L_g$  and  $D_\Phi(\mathbf{x}^* \|\mathbf{x}^1)$  is bounded by  $G^2$  then

$$\mathbb{E} \left[ \sum_{t=1}^T F^t(\mathbf{x}^t) \right] \geq \frac{\lambda}{\mu} \sum_{t=1}^T F^t(\mathbf{x}^*) - \frac{1}{\mu \cdot \eta} G^2 - \frac{\eta}{\mu \cdot 2\alpha_\Phi} T L_g^2$$

Choosing  $\eta = \frac{G}{L_g} \sqrt{\frac{2\alpha_\Phi}{T}}$ , the algorithm is  $(\frac{\lambda}{\mu}, R(T))$ -regret where  $R(T) = O(\frac{G L_g}{\mu} \sqrt{2\alpha_\Phi T})$ .  $\square$

### 3.2 Bandit Algorithm

In this section, we consider the bandit setting in which at every time  $t$  one can observe only the reward  $f^t(\mathbf{x}^t)$  where  $f^t$  is a bounded function defined on the convex set  $\mathcal{K} = [0, 1]^n$ . W.l.o.g., assume that  $f^t : [0, 1]^n \rightarrow [0, 1]$ . In our algorithm, we will consider a discretization of  $[0, 1]^n$  and the multilinear relaxations of functions  $f^t$  on these discrete points constructed as follows.

**Discretization and Multilinear Extension.** Let  $f : [0, 1]^n \rightarrow [0, 1]$  be a function. Consider a lattice  $\mathcal{L} = \{0, 2^{-M}, 2 \cdot 2^{-M}, \dots, \ell \cdot 2^{-M}, \dots, 1\}^n$  where  $0 \leq \ell \leq 2^M$  for some large parameter  $M$  as a discretization of  $[0, 1]^n$ .  $M$  is a constant parameter to be chosen later. Note that each  $x_i \in \{0, 2^{-M}, 2 \cdot 2^{-M}, \dots, \ell \cdot 2^{-M}, \dots, 1\}$  can be uniquely decomposed as  $x_i = \sum_{j=0}^M 2^{-j} y_{ij}$  where  $y_{ij} \in \{0, 1\}$ . By this observation, we lift the set  $[0, 1]^n \cap \mathcal{L}$  to the  $(n \times (M+1))$ -dim space. Specifically, define a bijective *lifting* map  $\text{lift} : [0, 1]^n \cap \mathcal{L} \rightarrow \{0, 1\}^{n \times (M+1)}$  such that each point  $(x_1, \dots, x_n) \in \mathcal{K} \cap \mathcal{L}$  is mapped to the unique point  $(y_{10}, \dots, y_{1M}, \dots, y_{n0}, \dots, y_{nM}) \in \{0, 1\}^{n \times (M+1)}$  where  $x_i = \sum_{j=0}^M 2^{-j} y_{ij}$ . Define function  $\tilde{f} : \{0, 1\}^{n \times (M+1)} \rightarrow [0, 1]$  such that  $\tilde{f}(\mathbf{1}_S) := f(\text{lift}^{-1}(\mathbf{1}_S))$ ; in other words,  $\tilde{f}(\mathbf{1}_S) = f(\mathbf{x})$  where  $\mathbf{x} \in [0, 1]^n \cap \mathcal{L}$  and  $\mathbf{1}_S = \text{lift}(\mathbf{x})$ . Note that  $\mathbf{1}_S$  with  $S \subset [n \times (M+1)]$  is a  $(n \times (M+1))$ -dim vector with  $(ij)^{\text{th}}$ -coordinate equal to 1 if  $(i, j) \in S$  and equal to 0 otherwise. Consider a multilinear extension  $F : [0, 1]^{n \times (M+1)} \rightarrow [0, 1]$  of  $\tilde{f}$  defined as follows.

$$F(\mathbf{z}) := \sum_{S \subset [n \times (M+1)]} \tilde{f}(\mathbf{1}_S) \prod_{(i,j) \in S} z_{ij} \prod_{(i,j) \notin S} (1 - z_{ij}).$$

By the definition,  $F(\mathbf{z})$  can be seen as  $\mathbb{E}[\tilde{f}(\mathbf{1}_S)]$  where the  $(ij)^{\text{th}}$ -coordinate of  $\mathbf{1}_S$  equals 1 (i.e.,  $(\mathbf{1}_S)_{ij} = 1$ ) with probability  $z_{ij}$ .

**Algorithm description.** Our algorithm, formally given in Algorithm 1, is inspired by algorithm SCRIBBLE [1] which has been derived in the context of bandit linear optimization. It has been observed that the gradient estimates of the functions in SCRIBBLE are unbiased only if those functions are linear; and that represents a main obstacle in order to derive an algorithm with optimal regret guarantee  $R(T) = \tilde{O}(\sqrt{T})$ . While aiming for the regret of  $\tilde{O}(\sqrt{T})$ , in our algorithm, we overcome this obstacle by considering at every step the gradient estimate of the multilinear extension of the reward function (construction above). That gradient estimate will be indeed proved to be unbiased. Incorporating that gradient estimate to the scheme in [1] and following our approach, we prove the regret guarantee of the algorithm.

---

**Algorithm 1** Algorithm in the bandit setting.

---

- 1: Let  $\Phi$  be a  $\nu$ -self-concordant function over  $[0, 1]^{n \times (M+1)}$ .
- 2: Let  $\mathbf{z}^1 \in \text{int}([0, 1]^{n \times (M+1)})$  such that  $\nabla \Phi(\mathbf{z}^1) = 0$ .
- 3: **for**  $t = 1$  to  $T$  **do**
- 4:   Let  $\mathbf{A}^t = [\nabla^2 \Phi(\mathbf{z}^t)]^{-1/2}$ .
- 5:   Pick  $\mathbf{u}^t \in \mathbb{S}_n$  uniformly random and set  $\mathbf{y}^t = \mathbf{z}^t + \mathbf{A}^t \mathbf{u}^t$ .
- 6:   Round  $\mathbf{y}^t$  to a random point  $\mathbf{1}_{S^t} \in \{0, 1\}^{n \times (M+1)}$  such that element  $(i, j)$  appears in  $S^t$  with probability  $y_{ij}^t$ .
- 7:   Play  $\mathbf{x}^t = \text{lifft}^{-1}(\mathbf{1}_{S^t})$  and receive the reward of  $f^t(\mathbf{x}^t)$ .
- 8:   Let  $\mathbf{g}^t = -n(M+1)f^t(\mathbf{x}^t)(\mathbf{A}^t)^{-1}\mathbf{u}^t$  and compute the solution  $\mathbf{z}^{t+1} \in [0, 1]^{n \times (M+1)}$  by applying the mirror descent framework on  $F^t$  (Section 3.1). Specifically,

$$\mathbf{z}^{t+1} = \arg \max_{\mathbf{z} \in [0, 1]^{n \times (M+1)}} \{ \langle \eta \mathbf{g}^t, \mathbf{z} - \mathbf{z}^t \rangle - D_\Phi(\mathbf{z} \parallel \mathbf{z}^t) \}.$$


---

**Analysis.** The remaining of the section is devoted to the analysis of Algorithm 1. For simplicity, until the end of this section, denote  $m = n(M+1)$ . Let  $\mathbb{B}_m$  and  $\mathbb{S}_m$  be the unit ball and the unit sphere in  $m$  dimensions, respectively. For a constant  $\delta$ , define  $\hat{F}_\delta(\mathbf{z}) := \mathbb{E}_{\mathbf{w} \in \mathbb{B}_m} [F(\mathbf{z} + \delta \mathbf{w})]$  where  $\mathbf{w}$  is drawn from a uniform distribution over  $\mathbb{B}_m$ . We first prove some technical lemmas by exploiting properties of multilinear extensions. These lemmas are similar to those needed to prove the regret guarantee of SCRIBBLE [1] but have some subtle differences because the functions we are considering are not linear (which is the case in [1]).

**Lemma 3** *It holds that  $\hat{F}_\delta(\mathbf{z}) = F(\mathbf{z})$ . Similarly, it also holds that  $\mathbb{E}_{\mathbf{u} \in \mathbb{S}_m} [F(\mathbf{z} + \delta \mathbf{u})] = F(\mathbf{z})$ .*

*Proof* Intuitively, the lemma holds since  $F(\mathbf{z})$  is linear w.r.t  $z_i$  for every  $i$ . In the following, we prove the first identity  $\hat{F}_\delta(\mathbf{z}) = F(\mathbf{z})$ . The second identity can be proved by exactly the same argument (using sphere instead of ball).

Consider a monomial  $g_k(\mathbf{z}) = z_1 z_2 \dots z_k$  for  $1 \leq k \leq m$ . We first prove by induction the claim that  $\mathbb{E}_{\mathbf{w} \in \mathbb{B}_m(r)} [g_k(\mathbf{z} + \delta \mathbf{w})] = g_k(\mathbf{z})$  for every ball  $\mathbb{B}_m(r)$  with radius  $r$  in dimension  $m$ . (By this notation,  $\mathbb{B}_m = \mathbb{B}_m(1)$ .) The base case  $k = 0$  is trivial. Assume that the induction hypothesis holds for  $g_{k-1}(\mathbf{z})$ . For any vector  $\mathbf{w} \in \mathbb{B}_m(r)$ , vector  $\mathbf{w}' = (-v_k, \mathbf{w}_{-k})$  is also in  $\mathbb{B}_m(r)$  and

$$\begin{aligned} g(\mathbf{z} + \delta \mathbf{w}) + g(\mathbf{z} + \delta \mathbf{w}') &= (z_k + \delta v_k) \cdot g_{k-1}(\mathbf{z}_{-k} + \delta \mathbf{w}_{-k}) + (z_k - \delta v_k) \cdot g_{k-1}(\mathbf{z}_{-k} + \delta \mathbf{w}_{-k}) \\ &= 2z_k \cdot g_{k-1}(\mathbf{z}_{-k} + \delta \mathbf{w}_{-k}) \end{aligned}$$

Note that, for a given  $|w_k|$ , uniformly random vectors  $(\pm w_k, \mathbf{w}_{-k})$  in the ball  $\mathbb{B}_m(r)$  induce uniformly random vectors  $\mathbf{w}_{-k}$  in the ball  $\mathbb{B}_{m-1}(\sqrt{r^2 - |w_k|^2})$ . Therefore,

$$\begin{aligned} \mathbb{E}_{\mathbf{w} \in \mathbb{B}_m(r)} [g_k(\mathbf{z} + \delta \mathbf{w})] &= z_k \cdot \mathbb{E}_{v_k} \mathbb{E}_{\mathbf{w}_{-k} \in \mathbb{B}_{m-1}(\sqrt{r^2 - |v_k|^2})} [g_{k-1}(\mathbf{z}_{-k} + \delta \mathbf{w}_{-k})] \\ &= z_k \cdot \mathbb{E}_{v_k} [g_{k-1}(\mathbf{z}_{-k})] = z_k \cdot g_{k-1}(\mathbf{z}_{-k}) = g_k(\mathbf{z}) \end{aligned}$$

where the second equality is due to the induction hypothesis. The claim follows.

As the multilinear extension is the sum of monomials multiplying with constant factors, the lemma holds because of the linearity of expectation.  $\square$

We restate here an useful lemma in [23].

**Lemma 4 ([23], Lemma 6.4)** *Let  $\delta > 0$  be a fixed constant and  $A \in \mathbb{R}^{m \times m}$  be an invertible matrix. Let  $G(\mathbf{z}) := F(A\mathbf{z})$  and  $\hat{G}_\delta(\mathbf{z}) := \mathbb{E}_{\mathbf{w} \in \mathbb{B}_m} [G(\mathbf{z} + \delta \mathbf{w})]$ . Then, it holds that*

$$\mathbb{E}_{\mathbf{u} \in \mathbb{S}_m} \left[ G(\mathbf{z} + \delta \mathbf{u}) \mathbf{u} \right] = \frac{\delta}{m} \nabla \hat{G}_\delta(\mathbf{z})$$

where the expectation is taken over uniform vector  $\mathbf{u}$  in the  $m$ -dim unit sphere  $\mathbb{S}_m$ .

**Lemma 5** Let  $A \in \mathbb{R}^{m \times m}$  be an invertible matrix. Define  $\hat{F}(z) := \mathbb{E}_{\mathbf{w} \in \mathbb{B}_m} [F(z + A\mathbf{w})]$ . Then it holds that

$$(i) \quad \hat{F}(z) = F(z) = \mathbb{E}_{\mathbf{u} \in \mathbb{S}_m} [F(z + A\mathbf{u})].$$

$$(ii) \quad \nabla \hat{F}(z) = m \mathbb{E}_{\mathbf{u} \in \mathbb{S}_m} [F(z + A\mathbf{u}) A^{-1} \mathbf{u}].$$

*Proof* We prove the first part of the lemma. Again, we prove only the identity  $F(z) = \mathbb{E}_{\mathbf{w} \in \mathbb{B}_m} [F(z + A\mathbf{w})]$ ; the identity  $F(z) = \mathbb{E}_{\mathbf{u} \in \mathbb{S}_m} [F(z + A\mathbf{u})]$  can be proved using exactly the same arguments. Define  $G(z) := F(Az)$ . Note that by this definition,  $F(z + A\mathbf{w}) = G(A^{-1}z + \mathbf{w})$ . The multilinear extension  $F$  is the weighted sum of monomials, so is  $G$ . Therefore, by the same argument as in the proof of Lemma 3, for any  $z$  and  $\delta$ , it holds that  $G(z) = \mathbb{E}_{\mathbf{w} \in \mathbb{B}_m} [G(z + \delta\mathbf{w})]$ . Applying this identity with  $\delta = 1$ , we have

$$G(A^{-1}z) = \mathbb{E}_{\mathbf{w} \in \mathbb{B}_m} [G(A^{-1}z + \mathbf{w})] \quad \Leftrightarrow \quad F(z) = \mathbb{E}_{\mathbf{w} \in \mathbb{B}_m} [F(z + A\mathbf{w})] = \hat{F}(z).$$

In the sequel, we prove the second part of the lemma. In fact, it can be proved using the same analysis in [23, Corollary 6.5]; we present it here for completeness. Define  $\hat{G}(z) := \mathbb{E}_{\mathbf{w} \in \mathbb{B}_m} [G(z + \mathbf{w})]$ . We have

$$\begin{aligned} m \mathbb{E}_{\mathbf{u} \in \mathbb{S}_m} [F(z + A\mathbf{u}) A^{-1} \mathbf{u}] &= mA^{-1} \mathbb{E}_{\mathbf{u} \in \mathbb{S}_m} [F(z + A\mathbf{u}) \mathbf{u}] = mA^{-1} \mathbb{E}_{\mathbf{u} \in \mathbb{S}_m} [G(A^{-1}z + \mathbf{u}) \mathbf{u}] \\ &= A^{-1} \nabla \hat{G}(A^{-1}z) = A^{-1} A \nabla \hat{F}(z) = \nabla \hat{F}(z) \end{aligned}$$

where the third inequality is due to Lemma 4 with  $\delta = 1$ .  $\square$

**Lemma 6** If  $f(x) \leq 1$  for every  $x \in [0, 1]^n$  then the corresponding multilinear extension  $F$  is  $2\sqrt{m}$ -Lipschitz.

*Proof* The proof comes directly from inspecting the derivatives. For each  $1 \leq \ell \leq m$ , we have

$$\begin{aligned} \left| \frac{\partial F(z)}{\partial z_\ell} \right| &= \left| \sum_{S: \ell \in S} \tilde{f}(\mathbf{1}_S) \prod_{j \in S} z_j \prod_{j \notin S} (1 - z_j) - \sum_{S: \ell \notin S} \tilde{f}(\mathbf{1}_S) \prod_{j \in S} z_j \prod_{j \notin S} (1 - z_j) \right| \\ &\leq \sum_{S: \ell \in S} \prod_{j \in S} z_j \prod_{j \notin S} (1 - z_j) + \sum_{S: \ell \notin S} \prod_{j \in S} z_j \prod_{j \notin S} (1 - z_j) \leq 2 \end{aligned}$$

where the first inequality is due to  $\tilde{f}(\mathbf{1}_S) \leq 1$  for all  $S \subset [m]$ . Therefore,  $\|\nabla F(z)\| \leq 2\sqrt{m}$ .  $\square$

**Theorem 2** Let  $f^t : [0, 1]^n \rightarrow [0, 1]$  be the reward function at time  $1 \leq t \leq T$  and let  $F^t$  be the multilinear extension of the discretization of  $f^t$  based on a lattice  $\mathcal{L}$  (defined earlier). Assume that  $F^t$ 's are  $(\lambda, \mu)$ -concave and for every  $x \in [0, 1]^n$ , there exists  $\bar{x} \in \mathcal{L}$  such that  $|f^t(x) - f^t(\bar{x})| \leq L \cdot 2^{-M}$  for every  $1 \leq t \leq T$  (for example,  $f^t$ 's are  $L$ -Lipschitz). Then, by choosing  $M = \log T$  and  $\eta = O\left(\frac{1}{(nM)^{3/2} \sqrt{T}}\right)$  and  $\Phi$  as a  $O(nM)$ -self-concordant function, Algorithm 1 achieves:

$$\sum_{t=1}^T \mathbb{E}[f^t(\mathbf{x}^t)] \geq \frac{\lambda}{\mu} \cdot \max_{\mathbf{x} \in [0, 1]^n} \sum_{t=1}^T f^t(\mathbf{x}) - O(\max\{\lambda/\mu, 1\} L n^{3/2} (\log T)^{3/2} (\log \log T) \sqrt{T}).$$

*Proof* A crucial point in the analysis is the observation that the gradient estimator of the multilinear relaxation is unbiased. Specifically,

$$\mathbb{E}[\mathbf{g}^t] = \mathbb{E}_{\mathbf{u}^t} \mathbb{E}_{\mathbf{x}^t} [\mathbf{g}^t] = \mathbb{E}_{\mathbf{u}^t} [m F^t(\mathbf{y}^t) (\mathbf{A}^t)^{-1} \mathbf{u}^t] = \nabla \hat{F}^t(\mathbf{z}^t) = \nabla F^t(\mathbf{z}^t)$$

where the second equality holds since  $\mathbb{E}[f^t(\mathbf{x}^t)] = F^t(\mathbf{y}^t)$  (by independent rounding and  $F^t$  is multilinear relaxation of  $f^t$ ); and the third and last equalities follow Lemma 5.

The remaining of the proof is similar to that in [23, Chapter 6] with some subtle differences because we consider the multilinear extensions of reward functions. Let  $0 < \delta < 1/2$  be some small constant and consider the hypercube  $[\delta, 1 - \delta]^m$ . Note that,  $[\delta, 1 - \delta]^m$  is convex and all balls of radius  $\delta$  around points in  $[\delta, 1 - \delta]^m$  are included in  $[0, 1]^m$ .

Let  $\mathbf{z}_\delta^*$  be the projection of  $\mathbf{z}^*$  onto  $[\delta, 1-\delta]^m$ . Then by properties of projections, we have  $\|\mathbf{z}_\delta^* - \mathbf{z}^*\| \leq \delta\sqrt{m}$ . Moreover, as  $|f^t(\mathbf{x})| \leq 1$  for every  $\mathbf{x}, t$ , and by definitions of the local norm, of  $\mathbf{A}^t$  and  $\mathbf{u}^t \in \mathbb{S}_m$ , it holds that

$$\|\mathbf{g}^t\|_{*\mathbf{x}^t}^2 \leq m^2 f^t(\mathbf{x}^t)^2 (\mathbf{u}^t)^\top ((\mathbf{A}^t)^{-1})^\top \nabla^{-2} \Phi(\mathbf{x}^t) (\mathbf{A}^t)^{-1} \mathbf{u}^t \leq m^2 \quad (6)$$

Besides, by Lemma 6, the multilinear relaxations  $F^t$ 's are  $(2\sqrt{m})$ -Lipschitz for every  $1 \leq t \leq T$ .

For any  $\mathbf{z}^* \in [0, 1]^m$ , we have

$$\begin{aligned} & \frac{\lambda}{\mu} \sum_{t=1}^T F^t(\mathbf{z}^*) - \sum_{t=1}^T \mathbb{E}_{\mathbf{g}^t, \mathbf{u}^t} [F^t(\mathbf{y}^t)] \\ &= \frac{\lambda}{\mu} \sum_{t=1}^T F^t(\mathbf{z}^*) - \sum_{t=1}^T \mathbb{E}_{\mathbf{g}^t} [F^t(\mathbf{z}^t)] && \text{by Lemma 5(i)} \\ &\leq \frac{\lambda}{\mu} \sum_{t=1}^T F^t(\mathbf{z}_\delta^*) - \sum_{t=1}^T \mathbb{E}_{\mathbf{g}^t} [F^t(\mathbf{z}^t)] + T(2\sqrt{m})\|\mathbf{z}^* - \mathbf{z}_\delta^*\| && 2\sqrt{m}\text{-Lipschitz of } F^t \\ &\leq \frac{1}{\mu \cdot \eta} D_\Phi(\mathbf{z}_\delta^* \|\mathbf{z}^1) + \frac{\eta}{\mu \cdot 2\alpha_\Phi} \sum_{t=1}^T \|\mathbf{g}^t\|_*^2 + 2\delta mT && \text{by Theorem 1} \\ &= \frac{1}{\mu \cdot \eta} [\Phi(\mathbf{z}_\delta^*) - \Phi(\mathbf{z}^1) - \langle \nabla \Phi(\mathbf{z}^1), \mathbf{z}_\delta^* - \mathbf{z}^1 \rangle] + \frac{\eta}{\mu \cdot 2\alpha_\Phi} m^2 T + 2\delta mT \\ &\leq \frac{1}{\mu \cdot \eta} [\Phi(\mathbf{z}_\delta^*) - \Phi(\mathbf{z}^1)] + \frac{\eta}{\mu \cdot 2\alpha_\Phi} m^2 T + 2\delta mT && \text{since } \nabla \Phi(\mathbf{z}^1) = 0 \\ &\leq \frac{\nu}{\mu \cdot \eta} \log \frac{1}{1 - \pi_{\mathbf{z}^1}(\mathbf{z}_\delta^*)} + \frac{\eta}{\mu \cdot 2\alpha_\Phi} m^2 T + 2\delta mT && \text{by Lemma 2} \\ &\leq \frac{\nu \log \frac{1}{\delta}}{\mu \cdot \eta} + \frac{\eta}{\mu \cdot 2\alpha_\Phi} m^2 T + 2\delta mT && \text{since } \mathbf{z}_\delta^* \in [\delta, 1 - \delta]^m \\ &= O(m\sqrt{\nu T} \log(mT)) \end{aligned}$$

The second inequality holds since the unbiased stochastic gradients  $\mathbf{g}^t$ 's have corresponding dual local norm bounded by  $m^2$ . The last equality follows the choice  $\eta = O(1/m\sqrt{\nu T})$ ,  $\delta = O(1/m\sqrt{T})$ .

Besides, by the property of multilinear extension and  $\mathbf{x}^t$  is obtained from  $\mathbf{y}^t$  by independently rounding,  $F^t(\mathbf{y}^t) = \mathbb{E}_{\mathbf{x}^t} [f^t(\mathbf{x}^t)]$ . Moreover,  $F^t$  and  $f^t$  have the same value on  $\{0, 1\}^n \cap \mathcal{L}$ . Therefore,

$$\begin{aligned} & \frac{\lambda}{\mu} \cdot \max_{\mathbf{x} \in [0, 1]^n \cap \mathcal{L}} \sum_{t=1}^T f^t(\mathbf{x}) - \sum_{t=1}^T \mathbb{E}_{\mathbf{x}^t, \mathbf{u}^t, \mathbf{g}^t} [f^t(\mathbf{x}^t)] = \frac{\lambda}{\mu} \cdot \max_{\mathbf{z} \in [0, 1]^m} \sum_{t=1}^T F^t(\mathbf{z}) - \sum_{t=1}^T \mathbb{E}_{\mathbf{u}^t} [F^t(\mathbf{y}^t)] \\ & \leq O(m\sqrt{\nu T} \log(mT)) = O(m^{3/2} \sqrt{T} \log(mT)) \end{aligned}$$

where the last equality is due to the fact that  $O(m)$ -self-concordant barrier exists [34, 9].

Let  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in [0, 1]^n} \sum_{t=1}^T f^t(\mathbf{x})$ . By the theorem assumption, there exists  $\bar{\mathbf{x}}^* \in [0, 1]^n \cap \mathcal{L}$  such that  $|f^t(\mathbf{x}) - f^t(\bar{\mathbf{x}})| \leq L \cdot 2^{-M}$ . Hence,  $\sum_{t=1}^T f^t(\mathbf{x}^*) \leq \sum_{t=1}^T f^t(\bar{\mathbf{x}}^*) + TL2^{-M}$ . Choose  $M = \log T$  (and recall that  $m = n(M + 1)$ ), we obtain the guarantee

$$\begin{aligned} & \frac{\lambda}{\mu} \cdot \max_{\mathbf{x} \in [0, 1]^n} \sum_{t=1}^T f^t(\mathbf{x}) - \sum_{t=1}^T \mathbb{E}_{\mathbf{x}^t, \mathbf{u}^t, \mathbf{g}^t} [f^t(\mathbf{x}^t)] \leq O\left(n^{3/2} M^{3/2} \sqrt{T} \log(nMT) + \frac{\lambda}{\mu} TL \frac{1}{2^M}\right) \\ & = O(\max\{\lambda/\mu, 1\} Ln^{3/2} (\log T)^{3/2} (\log \log T) \sqrt{T}). \end{aligned}$$

□

## 4 Perturbed Discrete Time Fictitious Play in Smooth Auctions

We consider adaptive dynamics in auctions. In the setting, there is an underlying auction  $\mathcal{o}$  and there are  $n$  players, each player  $i$  has a set of actions  $\mathcal{A}_i$  (that can be arbitrarily large but finite) and

a valuation function  $v_i$  taking values in  $[0, 1]$ . In each time step  $1 \leq t \leq T$ , each player  $i$  selects a strategy which is a distribution in the space of distributions  $\Delta(\mathcal{A}_i)$  according to some adaptive dynamic. The strategy profile at time  $t$  is denoted as  $\sigma^t \in \Delta(\mathcal{A})$ . Given the strategy profile  $\sigma^t$ , the auction induces a social welfare  $\text{Sw}(\sigma, \sigma^t) := \mathbb{E}_{\mathbf{a} \sim \sigma^t} [\text{Sw}(\sigma(\mathbf{a}); \mathbf{v})]$ . In this setting, we study the performance of adaptive dynamics, especially the ones which are not guaranteed to fulfill the vanishing regret condition, and eventually design dynamics/auctions with performance guarantee. Among others, *fictitious play* is an interesting, widely-studied dynamic which attracts a lot of attention in the community. In this section, we will study the performance of a version of fictitious play in smooth auctions.

**Valuation-Oriented Fictitious Play.** Consider the Perturbed Discrete Time Fictitious Play (PDTFP) — a smooth version of Discrete Time Fictitious Play (for example, see [26]). Let  $\Phi_i : \Delta(\mathcal{A}_i) \rightarrow \mathbb{R}$  for  $1 \leq i \leq n$  be strongly convex functions ( $\Phi_i$ 's are not necessarily the same). Initially, each player chooses some arbitrary action. At time  $t + 1$ , given a strategy profile  $\sigma^t$  where  $\sigma_i^t \in \Delta(\mathcal{A}_i)$  and perturbations  $N_i^t : \Delta(\mathcal{A}_i) \rightarrow \mathbb{R}^+$  for  $1 \leq i \leq n$  defined as  $N_i^t(\sigma_i) = D_{\Phi_i}(\sigma_i \parallel \sigma_i^t)$ , player  $i$  selects a mixed strategy  $\sigma_i^{t+1}$  such that

$$\sigma_i^{t+1} \in \arg \max_{\sigma_i \in \Delta(\mathcal{A}_i)} \mathbb{E}_{\mathbf{a}_i \sim \sigma_i} \mathbb{E}_{\mathbf{a}_{-i}^t \sim \sigma_{-i}^t} [v_i(\mathbf{a}_{-i}^t, a_i)] - \frac{1}{\eta} N_i^t(\sigma_i)$$

Equivalently,

$$\sigma_i^{t+1} \in \arg \max_{\sigma_i \in \Delta(\mathcal{A}_i)} \mathbb{E}_{\mathbf{a}_i \sim \sigma_i} \mathbb{E}_{\mathbf{a}^t \sim \sigma^t} [v_i(\mathbf{a}_{-i}^t, a_i) - v_i(\mathbf{a}^t)] - \frac{1}{\eta} N_i^t(\sigma_i), \quad (7)$$

since  $\mathbb{E}[v_i(\mathbf{a}^t)]$  is already determined. One common example of perturbations is the *relative entropy* (or *Kullback-Leibler divergence*), defined as

$$N_i^t(\sigma_i) = \sum_{a \in \mathcal{A}_i} \sigma_i(a) \log \frac{\sigma_i(a)}{\sigma_i^t(a)}.$$

which is the Bregman divergence with the negative entropy function  $\Phi_i(\sigma_i) = \sum_{a \in \mathcal{A}_i} \sigma_i(a) \log \sigma_i(a)$ .

Let  $V_i$  be the multilinear extension of the valuation  $v_i$  of player  $i$  (construction in Section 3.2) where now the corresponding lattice is the set of pure strategies  $\mathcal{A}$ . Note that the social welfare is the sum of all player valuations. Given an action profile  $\mathbf{a}^t$ , define  $\nabla^t(\mathbf{a}^t) : \mathbb{R}^n \rightarrow \mathbb{R}$  such as

$$\langle \nabla^t(\mathbf{a}^t), \mathbf{x} \rangle = \sum_{i=1}^n \frac{\partial V_i(\mathbf{a})}{\partial a_i} \cdot x_i.$$

As  $V_i$  is the multilinear extension of  $v_i$ , for every action  $\mathbf{a}^*$  we have

$$\langle \nabla^t(\mathbf{a}^t), \mathbf{a}^* - \mathbf{a}^t \rangle = \sum_{i=1}^n [V_i(a_i^*, \mathbf{a}_{-i}^t) - V_i(a_i^t, \mathbf{a}_{-i}^t)].$$

The PDTFP dynamic can be cast as the mirror descent algorithm. By Equation (7) — the update rules of PDTFP dynamic — at every time step  $t$ , strategy profile  $\sigma^{t+1} = (\sigma_1^{t+1}, \dots, \sigma_n^{t+1})$  is exactly the solution of the mirror descent update Equation (3):

$$\sigma^{t+1} \in \arg \max_{\sigma \in \Delta(\mathcal{A})} \mathbb{E}_{\mathbf{a} \sim \sigma} \mathbb{E}_{\mathbf{a}^t \sim \sigma^t} [\langle \nabla^t(\mathbf{a}^t), \mathbf{a} - \mathbf{a}^t \rangle] - \frac{1}{\eta} D_{\Phi}(\sigma \parallel \sigma^t)$$

where  $\Phi$  is a strongly convex function such that  $\Phi(\sigma) = \sum_{i=1}^n \Phi_i(\sigma_i)$ . Remark that, by the definition of  $\Phi$ ,  $D_{\Phi}(\sigma \parallel \sigma^t) = \sum_{i=1}^n D_{\Phi_i}(\sigma_i \parallel \sigma_i^t) = \sum_{i=1}^n N_i^t(\sigma_i)$ . Again, if  $\Phi_i$ 's are negative entropy functions (so the perturbations  $N_i^t$  are relative entropy functions) then  $\Phi(\sigma) = \sum_{i=1}^n \Phi_i(\sigma_i) = \sum_{i=1}^n \sum_{a \in \mathcal{A}_i} \sigma_i(a) \log \sigma_i(a)$ . Note that the PDTFP dynamic associated to that choice of entropy function is usually called *smooth fictitious play* [19] (or *logit dynamic*).

**PDTFP dynamic in smooth auctions.** Given arbitrary perturbations  $\Phi$ , it is not clear whether the sequences of player actions in PDTFP dynamic are individually-vanishing-regret, i.e., satisfying condition (2). In particular, in the dynamic players are valuation-oriented whereas the condition (2) concerns player utilities. Hence, the welfare guarantee by [36, 40] for individually-vanishing-regret sequences cannot be applied. In the following, we show that although it is not known whether PDTFP dynamic are individually-vanishing-regret, they achieve a guarantee bound on the welfare in smooth auctions.

**Theorem 3** *If the underlying auction  $\mathfrak{o}$  is a  $(\lambda, \mu)$ -smooth and  $D_{\Phi}(\cdot|\cdot)$  is bounded by  $G^2$  then the PDTFP dynamic with parameter  $\eta = O(G/\sqrt{T})$  achieves  $\left(\frac{\lambda}{1+\mu}, R(T)\right)$ -regret where  $R(T) = O\left(\frac{G\sqrt{T}}{1+\mu}\right)$ . In particular, if the perturbation is the relative entropy function then  $R(T) = O\left(\frac{\sqrt{T \log(n|\mathcal{A}|)}}{1+\mu}\right)$ .*

*Proof* The analysis follows closely the one of Theorem 1 with some modifications. For simplicity, without loss of generality, assume that the distributions  $\bar{D}_1, \dots, \bar{D}_n$  in the definition of smooth auctions (Definition 3) give rise to a pure strategy profile  $\bar{\mathbf{a}}$  and at any time step  $t$ , the PDTFP dynamic outputs a pure profile  $\mathbf{a}^t$ . The analysis remains the same for general distributions/mixed profiles by putting additional expectations into some formula.

As the underlying auction is  $(\lambda, \mu)$ -smooth, given a fixed valuation profile  $\mathbf{v}$ , there exists a strategy profile  $\bar{\mathbf{a}}$  such that for any profile  $\mathbf{a}$ , it holds that

$$\sum_{i=1}^n u_i(\bar{a}_i, \mathbf{a}_{-i}; v_i) \geq \lambda \cdot \text{OPT}(\mathbf{v}) - \mu \cdot \text{REV}(\mathbf{a})$$

where  $\text{OPT}(\mathbf{v})$  stands for the optimal welfare given the valuation profile  $\mathbf{v}$ . We first derive an useful inequality based on the smoothness of the auction. We have

$$\begin{aligned} \langle \nabla^t(\mathbf{a}^t), \bar{\mathbf{a}} - \mathbf{a}^t \rangle &= \sum_i [V_i(\bar{a}_i, \mathbf{a}_{-i}^t; \mathbf{v}) - V_i(a_i, \mathbf{a}_{-i}^t; \mathbf{v})] \\ &= \sum_i [u_i(\bar{a}_i, \mathbf{a}_{-i}^t; v_i) + p_i(\bar{a}_i, \mathbf{a}_{-i}^t; v_i)] - \text{SW}(\mathbf{a}^t; \mathbf{v}) \\ &\geq \lambda \cdot \text{OPT}(\mathbf{v}) - \mu \cdot \text{REV}(\mathbf{a}^t; \mathbf{v}) - \text{SW}(\mathbf{a}^t; \mathbf{v}) \\ &\geq \lambda \cdot \text{OPT}(\mathbf{v}) - (1 + \mu) \cdot \text{SW}(\mathbf{a}^t). \end{aligned} \tag{8}$$

The first inequality follows by the  $(\lambda, \mu)$ -smoothness and the non-negativity of payments  $p_i$ 's. The second inequality is obvious since the revenue is always smaller than the welfare. We remark that Inequality (8) is similar to (but not the same as) the notation of  $(\cdot, \cdot)$ -concavity since it can written as

$$\langle \nabla \text{SW}(\mathbf{a}^t), \bar{\mathbf{a}} - \mathbf{a}^t \rangle \geq \lambda \cdot \text{SW}(\mathbf{a}^*) - (1 + \mu) \cdot \text{SW}(\mathbf{a}^t)$$

where  $\mathbf{a}^*$  is the optimal strategy. Hence, there would be a connection between concavity and smoothness.

Define the potential as  $\Psi^t = \frac{1}{\eta} D_{\Phi}(\bar{\mathbf{a}}|\mathbf{a}^t)$ . Note that here we use the Bregman divergence from the strategy  $\bar{\mathbf{a}}$  (induced by the smooth auction) to  $\mathbf{a}^t$  instead of the Bregman divergence from the optimal strategy  $\mathbf{a}^*$  to  $\mathbf{a}^t$  (as in Theorem 1). By the same arguments proving Inequality (4), we have

$$\eta(\Psi^{t+1} - \Psi^t) = D_{\Phi}(\bar{\mathbf{a}}|\mathbf{a}^{t+1}) - D_{\Phi}(\bar{\mathbf{a}}|\mathbf{a}^t) \leq -\eta \langle \nabla^t(\mathbf{a}^t), \bar{\mathbf{a}} - \mathbf{a}^t \rangle + \frac{\eta^2}{2\alpha_{\Phi}} \|\nabla^t(\mathbf{a}^t)\|_*^2$$

Given the valuation profile  $\mathbf{v}$ , let  $\mathbf{a}^*$  be the action that gives the optimal welfare, i.e.,  $\text{Sw}(\mathbf{a}^*; \mathbf{v}) = \text{OPT}(\mathbf{v})$ . Using the same arguments as in the proof of Theorem 1, we have

$$\begin{aligned} \sum_{t=1}^T (\lambda \text{Sw}(\mathbf{a}^*) - (1 + \mu) \text{Sw}(\mathbf{a}^t)) &\leq \Psi^1 + \sum_{t=1}^T \left[ \lambda \text{Sw}(\mathbf{a}^*) - (1 + \mu) \text{Sw}(\mathbf{a}^t) + \Psi^{t+1} - \Psi^t \right] \\ &\leq \Psi^1 + \sum_{t=1}^T \left[ \underbrace{\lambda \text{OPT}(\mathbf{v}) - (1 + \mu) \text{Sw}(\mathbf{a}^t) - \langle \nabla^t(\mathbf{a}^t), \bar{\mathbf{a}} - \mathbf{a}^t \rangle}_{\leq 0 \text{ by Inequality (8)}} + \frac{\eta}{2\alpha_\Phi} \|\nabla_t(\mathbf{a}^t)\|_*^2 \right] \\ &\leq \frac{1}{\eta} D_\Phi(\bar{\mathbf{a}} \|\mathbf{a}^1) + \frac{\eta}{2\alpha_\Phi} \sum_{t=1}^T \|\nabla^t(\mathbf{a}^t)\|_*^2 \end{aligned}$$

Thus,

$$\sum_{t=1}^T \text{Sw}(\mathbf{a}^t) \geq \frac{\lambda}{1 + \mu} \sum_{t=1}^T \text{Sw}(\mathbf{a}^*) - \frac{1}{(1 + \mu)\eta} D_\Phi(\bar{\mathbf{a}} \|\mathbf{a}^1) - \frac{\eta}{(1 + \mu)2\alpha_\Phi} \sum_{t=1}^T \|\nabla^t(\mathbf{a}^t)\|_*^2$$

Note that if player valuations are in the range  $[0, 1]$ , then

$$\|\nabla^t(\mathbf{a}^t)\|_* \leq \|\nabla^t(\mathbf{a}^t)\|_\infty \leq 1.$$

By the theorem assumptions,  $D_\Phi(\bar{\mathbf{a}} \|\mathbf{a}^1) \leq G^2$ . Hence, choosing  $\eta = O(G/\sqrt{T})$ , the PDTFP dynamic achieves  $(\frac{\lambda}{1+\mu}, R(T))$ -regret where  $R(T) = O(\frac{G\sqrt{T}}{1+\mu})$ .

Consider the particular PDTFP dynamic with relative entropy perturbation. Function  $\Phi(\boldsymbol{\sigma})$  is  $\alpha_\Phi = \frac{1}{2 \ln 2}$ -strongly convex (due to Pinsker's inequality). Moreover,  $D_\Phi(\bar{\mathbf{a}} \|\mathbf{a}^1) \leq \max_i \log(n|\mathcal{A}_i|) \leq \log(n|\mathcal{A}|)$ . Therefore, choosing  $\eta = O(1/\sqrt{T \log(n|\mathcal{A}|)})$ , the PDTFP dynamic with relative entropy perturbation achieves  $(\frac{\lambda}{1+\mu}, R(T))$ -regret where  $R(T) = O(\frac{\sqrt{T \log(n|\mathcal{A}|)}}{1+\mu})$ .  $\square$

## 5 Online Simultaneous Second-Price Auctions with Reserve Prices

In this section, we are interested in the objective of maximizing the revenue. In the setting, there are  $n$  bidders and  $m$  items to be sold to these bidders. At each time step  $t = 1, 2, \dots, T$ , the auctioneer selects reserve prices  $r_i^t = (r_{i1}^t, \dots, r_{im}^t)$  for each bidder  $i$  where  $r_{ij}$  is the reserve price of item  $j$  for bidder  $i$ . Subsequently, every bidder  $i$  picks a bid vector  $\mathbf{b}_i^t = (b_{i1}^t, \dots, b_{im}^t)$  where  $b_{ij}^t$  is the bid of bidder  $i$  on item  $1 \leq j \leq m$ . Note that  $b_{ij}^t$  and  $b_{i'j}^t$  can be correlated. Then the auction for each item  $1 \leq j \leq m$  works as follows: (1) remove all bidders  $i$  with  $b_{ij}^t < r_{ij}^t$ ; (2) run the second-price auction on the remaining bidders to determine the winner of item  $j$ ; (3) charge the winner of item  $j$  the larger of  $r_{ij}^t$  and the second highest bid among the bids  $b_{ij}^t$  of remaining bidders. Denote the revenue of selling item  $j$  as  $\text{REV}_j(\mathbf{r}^t, \mathbf{b}^t)$  where  $\mathbf{b}^t = (b_1^t, \dots, b_n^t)$  and  $\mathbf{r}^t = (r_1^t, \dots, r_n^t)$ . The revenue of the auctioneer at time step  $t$  is  $\text{REV}(\mathbf{r}^t, \mathbf{b}^t) = \sum_{j=1}^m \text{REV}_j(\mathbf{r}^t, \mathbf{b}^t)$ . The goal of the auctioneer is to achieve the total revenue approximately close to that achieved by the best fixed reserve-price in hindsight  $\sum_{j=1}^m \text{REV}_j(\mathbf{r}^*, \mathbf{b}^t)$ .

In the setting, by scaling, assume that all bids and reserve prices are in  $\mathcal{K} = [0, 1]^{n \times m}$ . Consider the lattice  $\mathcal{L} = \{\ell \cdot 2^{-M} : 0 \leq \ell \leq 2^M\}^{n \times m} \subset [0, 1]^{n \times m}$  for some large parameter  $M$  as a discretization of  $[0, 1]^{n \times m}$ . Observe that for any reserve price vector  $\mathbf{r}$ ,  $|\text{REV}(\mathbf{r}, \mathbf{b}) - \text{REV}(\bar{\mathbf{r}}, \mathbf{b})| \leq m \cdot 2^{-M}$  where  $\bar{\mathbf{r}}$  is a reserve price vector such that  $\bar{r}_{ij}$  is the largest multiple of  $2^{-M}$  smaller than  $r_{ij}$  for every  $i, j$  (for some large enough  $M$ ). Therefore, one can approximate the revenue up to any arbitrary precision by restricting the reserve price on  $\mathcal{L}$ . We slightly abuse notation by denoting  $\text{REV}_j(\mathbf{1}_S, \mathbf{b})$  as  $\text{REV}_j(\mathbf{r}, \mathbf{b})$  where  $\mathbf{1}_S = m(\mathbf{r})$  (recall that  $m$  is the map defined in Section 3.2). Following Section 3.2, given a bid vector  $\mathbf{b}$ , the multilinear extension  $\overline{\text{REV}}$  of the revenue  $\text{REV}$  is defined as  $\overline{\text{REV}}(\cdot, \mathbf{b}) : [0, 1]^{n \times m \times (M+1)} \rightarrow \mathbb{R}$  such that:

$$\overline{\text{REV}}(\mathbf{z}, \mathbf{b}) = \sum_{S \subset [n \times m \times (M+1)]} \left( \sum_{j=1}^m \text{REV}_j(\mathbf{1}_S, \mathbf{b}) \right) \prod_{(i,j,k) \in S} z_{ijk} \prod_{(i,j,k) \notin S} (1 - z_{ijk}).$$

**Online bandit Reserve-Price Algorithm.** Initially, let  $\mathbf{r}^1$  be an arbitrary feasible reserve-price. At each time step  $t \geq 1$ ,

- (i) select  $\mathbf{r}^t$  or  $\mathbf{0}$  each with probability 1/2 as the reserve-price;
- (ii) receive the revenue corresponding to the selected reserve-price;
- (iii) compute  $\mathbf{r}^{t+1}$  using Algorithm 1 with the following specification: in line 8 of Algorithm 1, replace  $f^t(\mathbf{x}^t)$  by  $2\text{REV}(\mathbf{r}^t, \mathbf{b}^t)$  if the selected reserve-price is  $\mathbf{r}^t$ , or replace  $f^t(\mathbf{x}^t)$  by 0 if the selected reserve-price is  $\mathbf{0}$ . (By doing that, the expected value of  $\mathbf{g}^t$  in Algorithm 1 is  $-\nabla \overline{\text{REV}}(\mathbf{r}^t, \mathbf{b}^t)$ .)

**Analysis.** In order to analyze the performance of this algorithm, we study the properties of some related functions and then derive the regret bound for the algorithm.

Fix a bid vector  $\mathbf{b}$ . Let  $\mathbf{r}_j$  be a vector consisting of reserve prices on item  $j$ , i.e.,  $\mathbf{r}_j = (r_{1j}, \dots, r_{nj})$ . As  $\mathbf{b}$  is fixed and the selling procedure of each item depends only on the reserve prices to the item, so for simplicity denote  $\text{REV}_j(\mathbf{r}, \mathbf{b})$  as  $\text{REV}_j(\mathbf{r}_j)$  and  $\text{REV}(\mathbf{r}, \mathbf{b})$  as  $\text{REV}(\mathbf{r})$ . Define a function  $h_j : \{0, 1\}^{n \times (M+1)} \rightarrow \mathbb{R}$  such that  $h_j(\mathbf{1}_T) = \max\{\text{REV}_j(\mathbf{1}_T), \text{REV}_j(\mathbf{1}_\emptyset)\} = \max\{\text{REV}_j(\mathbf{r}), \text{REV}_j(\mathbf{0})\}$  where  $\mathbf{r}_j$  is the reserve price corresponding to  $\mathbf{1}_T$  for  $T \subset [n \times (M+1)]$ . Let  $H_j : [0, 1]^{n \times (M+1)} \rightarrow \mathbb{R}$  be the multilinear extension of  $h_j$ . Moreover, define  $H : [0, 1]^{n \times m \times (M+1)} \rightarrow \mathbb{R}$  as the multilinear extension of  $\max\{\text{REV}(\mathbf{r}), \text{REV}(\mathbf{0})\}$  defined as

$$H(\mathbf{z}) = \sum_{S \subset [n \times m \times (M+1)]} \max\{\text{REV}(\mathbf{1}_S), \text{REV}(\mathbf{1}_\emptyset)\} \prod_{(i,j,k) \in S} z_{ijk} \prod_{(i,j,k) \notin S} (1 - z_{ijk})$$

**Lemma 7** *It holds that  $H(\mathbf{z}) = \sum_{j=1}^m H_j(\mathbf{z}_j)$  where  $\mathbf{z}_j$  is the restriction of  $\mathbf{z}$  to the coordinate related to item  $j$ .*

*Proof* As items are sold separately,

$$H(\mathbf{z}) = \sum_{S \subset [n \times m \times (M+1)]} \left( \sum_{j=1}^m h_j(\mathbf{1}_A) \right) \prod_{(i,j,k) \in S} z_{ijk} \prod_{(i,j,k) \notin S} (1 - z_{ijk})$$

where  $A \subset [n \times (M+1)]$  is the restriction of  $S$  on coordinates related to item  $j$ . Therefore,

$$\begin{aligned} H(\mathbf{z}) &= \sum_{j=1}^m \sum_{U \subset [n \times (m-1) \times (M+1)]} \underbrace{\left[ \sum_{A \subset [n \times (M+1)]} h_j(\mathbf{1}_A) \prod_{(i,k) \in A} z_{ijk} \prod_{(i,k) \notin A} (1 - z_{ijk}) \right]}_{\text{independent of } U \text{ since the allocation of } j \text{ depends only on bids to item } j} \\ &\quad \cdot \prod_{(i,j',k) \in U} z_{ij'k} \prod_{(i,j',k) \notin U, j' \neq j} (1 - z_{ij'k}) \\ &= \sum_{j=1}^m \left[ \sum_{A \subset [n \times (M+1)]} h_j(\mathbf{1}_A) \prod_{(i,k) \in A} z_{ijk} \prod_{(i,k) \notin A} (1 - z_{ijk}) \right] \\ &\quad \cdot \underbrace{\sum_{U \subset [n \times (m-1) \times (M+1)]} \prod_{(i,j',k) \in U} z_{ij'k} \prod_{(i,j',k) \notin U, j' \neq j} (1 - z_{ij'k})}_{=1} \\ &= \sum_{j=1}^m \left[ \sum_{A \subset [n \times (M+1)]} h_j(\mathbf{1}_A) \prod_{(i,k) \in A} z_{ijk} \prod_{(i,k) \notin A} (1 - z_{ijk}) \right] = \sum_{j=1}^m H_j(\mathbf{z}_j) \end{aligned}$$

□

We will prove that  $H$  is  $(1, 1)$ -concave. By Lemma 7, it is sufficient to prove that property for every function  $H_j$ .

**Lemma 8** *Function  $H_j$  is  $(1, 1)$ -concave.*

*Proof* We prove that the condition (1) of the (1, 1)-concavity holds for all points in the lattice. As the multilinear extension can be seen as the expectation over these points, the lemma will follow. Fix a bid profile  $\mathbf{b}_j = (b_{1j}, \dots, b_{nj})$ . Without loss of generality, assume that  $b_{1j} \geq b_{2j} \geq \dots \geq b_{nj}$ . Let  $\mathbf{r}_j$  and  $\mathbf{r}_j^*$  be two arbitrary reserve price vectors. We will show that

$$\begin{aligned} \sum_{i=1}^n \left[ \max\{\text{REV}_j(\mathbf{r}_{-ij}, r_{ij}^*), \text{REV}_j(\mathbf{0})\} - \max\{\text{REV}_j(\mathbf{r}_j), \text{REV}_j(\mathbf{0})\} \right] \\ \geq \max\{\text{REV}_j(\mathbf{r}_j^*), \text{REV}_j(\mathbf{0})\} - \max\{\text{REV}_j(\mathbf{r}_j), \text{REV}_j(\mathbf{0})\} \end{aligned} \quad (9)$$

where  $\mathbf{r}_{-ij}$  stands for the reserve price vectors on item  $j$  without the reserve price of bidder  $i$ .

Observe that the revenue  $\max\{\text{REV}_j(\mathbf{r}'_j), \text{REV}_j(\mathbf{0})\}$  for every reserve price  $\mathbf{r}'_j$  is at least the second highest bid  $b_{2j}$  (that is obtained in  $\text{REV}_j(\mathbf{0})$ ). Moreover, for any reserve price  $\mathbf{r}'_j$  such that the auctioneer either (1) removes the first bidder (with highest bid) or (2) removes the second bidder and  $r'_{1j} \leq b_{2j}$ , the revenue

$$\max\{\text{REV}_j(\mathbf{r}'_j), \text{REV}_j(\mathbf{0})\} = \text{REV}_j(\mathbf{0}).$$

Hence,  $\max\{\text{REV}_j(\mathbf{r}'_j), \text{REV}_j(\mathbf{0})\} \neq \text{REV}_j(\mathbf{0})$  if and only if  $b_{2j} < r'_{1j} \leq b_{1j}$ .

By these observations, we deduce that

$$\max\{\text{REV}_j(\mathbf{r}_{-ij}, r_{ij}^*), \text{REV}_j(\mathbf{0})\} \neq \max\{\text{REV}_j(\mathbf{r}_j), \text{REV}_j(\mathbf{0})\}$$

if and only if  $i = 1$  and

- either  $b_{2j} \leq r_{1j} \neq r_{1j}^* \leq b_{1j}$ ;
- or  $r_{1j}^* \in (b_{2j}, b_{1j}]$  but  $r_{1j} \notin (b_{2j}, b_{1j}]$ ;
- or inversely  $r_{1j} \in (b_{2j}, b_{1j}]$  but  $r_{1j}^* \notin (b_{2j}, b_{1j}]$ .

Thus, proving Inequality (9) is equivalent to showing that

$$\begin{aligned} \max\{\text{REV}_j(\mathbf{r}_{-1j}, r_{1j}^*), \text{REV}_j(\mathbf{0})\} - \max\{\text{REV}_j(\mathbf{r}_j), \text{REV}_j(\mathbf{0})\} \\ \geq \max\{\text{REV}_j(\mathbf{r}_j^*), \text{REV}_j(\mathbf{0})\} - \max\{\text{REV}_j(\mathbf{r}_j), \text{REV}_j(\mathbf{0})\} \end{aligned}$$

**Case 1:**  $b_{2j} \leq r_{1j} \neq r_{1j}^* \leq b_{1j}$ . In this case, both sides are equal to  $r_{1j}^* - r_{1j}$ .

**Case 2:**  $r_{1j}^* \in (b_{2j}, b_{1j}]$  but  $r_{1j} \notin (b_{2j}, b_{1j}]$ . In this case, both sides are equal to  $r_{1j}^* - b_{2j}$ .

**Case 3:**  $r_{1j} \in (b_{2j}, b_{1j}]$  but  $r_{1j}^* \notin (b_{2j}, b_{1j}]$ . In this case, both sides are equal to  $b_{2j} - r_{1j}$ .

**Case 4:** the complementary of all previous cases. In this case, both sides are equal to 0.

Therefore, Inequality (9) holds and so the lemma follows.  $\square$

**Theorem 4** *The online bandit reserve price algorithm achieves the regret bound of  $(1/2, O(mn^{3/2}(\log T)^{3/2}(\log \log T)\sqrt{T}))$ .*

*Proof* Consider an imaginary algorithm which is similar to our online reserve price algorithm but at every step  $t$ , its gain on item  $j$  is  $\max\{\text{REV}_j(\mathbf{r}^t), \text{REV}_j(\mathbf{0})\}$ . (This algorithm is called imaginary since one cannot decide which reserve price between  $\mathbf{r}^t$  and  $\mathbf{0}$  is better when the bid vector is not known.) We verify the conditions of Theorem 2. The discretization satisfies the condition that for any given bids  $\mathbf{b}$ , for any reserve price  $\mathbf{r}$ , there exists a reserve price  $\bar{\mathbf{r}}$  in the lattice which gives  $|\text{REV}(\mathbf{r}, \mathbf{b}) - \text{REV}(\bar{\mathbf{r}}, \mathbf{b})| \leq m \cdot 2^{-M}$ . Moreover, Lemma 8 shows the (1, 1)-concavity of  $H$ . Therefore, applying Theorem 2, the imaginary algorithm achieves the regret bound of  $(1, R(T))$  where

$$R(T) = O(mn^{3/2}(\log T)^{3/2}(\log \log T)\sqrt{T}).$$

As the online reserve price algorithm selects at every step  $t$  either  $\mathbf{r}^t$  or  $\mathbf{0}$  with probability 1/2, the revenue of the algorithm is at least half that of the imaginary algorithm. The theorem follows.  $\square$

## 6 Bandit No-Envy Learning in Auctions

In this section we consider the bandit item selection problem. In the problem, there are  $m$  items and a player with monotone submodular valuation  $v : 2^{[m]} \rightarrow [0, 1]$ . At every time step  $1 \leq t \leq T$ , the player chooses a subset of items  $S^t \subset [m]$  and the adversary picks adaptively (probably depending on the history up to time  $t - 1$  but not on the current set  $S^t$ ) a threshold vector  $\mathbf{p}^t$ . The player observes only the thresholds  $p_j^t$  for  $j \in S^t$  and gets reward  $v(S^t) - \sum_{j \in S^t} p_j^t$ . Without loss of generality, assume that  $0 \leq v(S) \leq 1$  for all  $S \subset [m]$  and also  $0 \leq p_j^t \leq 1$  for all  $t$  and  $j$ .

We seek an  $r$ -approximate no-envy learning algorithm for some constant  $0 < r \leq 1$ . That is, for any adaptively chosen sequence of threshold vectors  $\mathbf{p}^t$  for  $1 \leq t \leq T$ , the sets  $S^t$  for  $1 \leq t \leq T$  chosen by the algorithm satisfy

$$\mathbb{E} \left[ \sum_{t=1}^T \left( v(S^t) - \sum_{j \in S^t} p_j^t \right) \right] \geq \max_{S \subseteq [m]} \sum_{t=1}^T \left( r \cdot v(S) - \sum_{j \in S} p_j^t \right) - R(T)$$

where the regret  $R(T) = o(T)$ .

Let  $V : [0, 1]^m \rightarrow \mathbb{R}^+$  be the multilinear relaxation of the monotone submodular valuation  $v$ . Formally,

$$V(\mathbf{z}) = \sum_{S \subseteq [m]} v(S) \prod_{j \in S} z_j \prod_{j \notin S} (1 - z_j).$$

It is well-known that if  $v$  is monotone submodular then  $V$  is also monotone and it satisfies the diminishing return property:  $\nabla V(\mathbf{x}) \geq \nabla V(\mathbf{y})$  for all  $\mathbf{x}, \mathbf{y} \in [0, 1]^m$  such that  $\mathbf{x} \leq \mathbf{y}$ . Here for two vectors  $\mathbf{a}, \mathbf{b}$ , we mean  $\mathbf{a} \leq \mathbf{b}$  iff  $a_j \leq b_j$  for all  $j$ .

The following lemma, which has been implicitly proved in [22], shows that the multilinear relaxation  $V$  is (1, 2)-concave.

**Lemma 9 ([22])** *For every  $\mathbf{x}, \mathbf{y} \in [0, 1]^m$ , it holds that*

$$\langle \nabla V(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \geq V(\mathbf{y}) - 2V(\mathbf{x})$$

*Proof* Given vector  $\mathbf{x}, \mathbf{y}$ , let  $\mathbf{x} \vee \mathbf{y}$  be the vector such that its  $i^{\text{th}}$  coordinate is  $\max\{x_i, y_i\}$  for every  $1 \leq i \leq n$ . Similarly, let  $\mathbf{x} \wedge \mathbf{y}$  be the vector such that its  $i^{\text{th}}$  coordinate is  $\min\{x_i, y_i\}$  for every  $1 \leq i \leq n$ .

For any vectors  $\mathbf{x} \leq \mathbf{z}$ , using the diminishing return property  $\nabla V(\mathbf{x}) \geq \nabla V(\mathbf{x} + t(\mathbf{z} - \mathbf{x}))$  for  $0 \leq t \leq 1$ , we have

$$\begin{aligned} V(\mathbf{z}) - V(\mathbf{x}) &= \int_0^1 \langle \mathbf{z} - \mathbf{x}, \nabla V(\mathbf{x} + t(\mathbf{z} - \mathbf{x})) \rangle dt \\ &\leq \int_0^1 \langle \mathbf{z} - \mathbf{x}, \nabla V(\mathbf{x}) \rangle dt = \langle \mathbf{z} - \mathbf{x}, \nabla V(\mathbf{x}) \rangle \end{aligned}$$

Therefore,

$$V(\mathbf{x} \vee \mathbf{y}) - V(\mathbf{x}) \leq \langle \mathbf{x} \vee \mathbf{y} - \mathbf{x}, \nabla V(\mathbf{x}) \rangle. \quad (10)$$

Similarly for vectors  $\mathbf{x} \leq \mathbf{z}$ , we have

$$\begin{aligned} V(\mathbf{z}) - V(\mathbf{x}) &= \int_0^1 \langle \mathbf{z} - \mathbf{x}, \nabla V(\mathbf{x} + t(\mathbf{z} - \mathbf{x})) \rangle dt \\ &\geq \int_0^1 \langle \mathbf{z} - \mathbf{x}, \nabla V(\mathbf{z}) \rangle dt = \langle \mathbf{z} - \mathbf{x}, \nabla V(\mathbf{z}) \rangle. \end{aligned}$$

Therefore,

$$V(\mathbf{x} \wedge \mathbf{y}) - V(\mathbf{x}) \leq \langle \mathbf{x} \wedge \mathbf{y} - \mathbf{x}, \nabla V(\mathbf{x}) \rangle \quad (11)$$

Summing (10) and (11) and using the fact  $(\mathbf{x} \vee \mathbf{y}) + (\mathbf{x} \wedge \mathbf{y}) = \mathbf{x} + \mathbf{y}$ , we obtain

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

As  $V$  is monotone and non-negative, we deduce that

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

□

In order to apply our framework, we first prove the guarantee of the mirror-descent algorithm similar to the one in Section 3.1.

**Mirror descent.** Let  $\Phi$  be a  $\alpha_\Phi$ -strongly convex function w.r.t  $\|\cdot\|$ . Initially, let  $\mathbf{z}^1$  is an arbitrary feasible point. At time step  $t$ , play  $\mathbf{z}^t$  and receive the vector  $\mathbf{p}^t$ . Compute  $-\mathbf{g}^t$  an unbiased estimate of  $\frac{1}{2}\nabla(V(\mathbf{z}^t) - \langle \mathbf{p}^t, \mathbf{z}^t \rangle) = \frac{1}{2}\nabla V(\mathbf{z}^t) - \mathbf{p}^t$ . Update  $\mathbf{z}^{t+1}$  as follows:

$$\mathbf{z}^{t+1} = \arg \max_{\mathbf{z} \in [0,1]^m} \{ \langle \eta \mathbf{g}^t, \mathbf{z} - \mathbf{z}^t \rangle - D_\Phi(\mathbf{z} \|\mathbf{z}^t) \}$$

**Theorem 5** Then the mirror descent algorithm above achieves

$$\sum_{t=1}^T \left( V(\mathbf{z}^t) - \langle \mathbf{p}^t, \mathbf{z}^t \rangle \right) \geq \max_{\mathbf{z} \in [0,1]^n} \sum_{t=1}^T \left( \frac{1}{2} V(\mathbf{z}) - \langle \mathbf{p}^t, \mathbf{z} \rangle \right) - \frac{1}{\eta} D_\Phi(\mathbf{x}^* \|\mathbf{x}^1) - \frac{\eta}{2\alpha_\Phi} \sum_{t=1}^T \|\mathbf{g}^t\|_*^2$$

*Proof* The analysis is similar to that of Theorem 1. Let  $\mathbf{z}^* \in \arg \max_{\mathbf{z} \in [0,1]^n} \sum_{t=1}^T \left( \frac{1}{2} V(\mathbf{z}) - \langle \mathbf{p}^t, \mathbf{z} \rangle \right)$ . Define the potential as  $\Psi^t = \frac{1}{\eta} D_\Phi(\mathbf{z}^* \|\mathbf{z}^t)$ . By the same argument as in the analysis of Theorem 1, we have

$$D_\Phi(\mathbf{z}^* \|\mathbf{z}^{t+1}) - D_\Phi(\mathbf{z}^* \|\mathbf{z}^t) \leq \frac{\eta^2}{2\alpha_\Phi} \|\mathbf{g}^t\|_*^2 - \eta \langle \mathbf{g}^t, \mathbf{z}^* - \mathbf{z}^t \rangle \quad (12)$$

Using the bound of the potential change due to Inequality (12), we get

$$\begin{aligned} & \sum_{t=1}^T \left( \frac{1}{2} V(\mathbf{z}^*) - \langle \mathbf{p}^t, \mathbf{z}^* \rangle - V(\mathbf{z}^t) + \langle \mathbf{p}^t, \mathbf{z}^t \rangle \right) \\ & \leq \Psi_1 + \sum_{t=1}^T \left( \frac{1}{2} V(\mathbf{z}^*) - V(\mathbf{z}^t) - \langle \mathbf{p}^t, \mathbf{z}^* - \mathbf{z}^t \rangle + \Psi^{t+1} - \Psi^t \right) \\ & \leq \Psi_1 + \sum_{t=1}^T \left( \frac{1}{2} V(\mathbf{z}^*) - V(\mathbf{z}^t) - \langle \mathbf{p}^t, \mathbf{z}^* - \mathbf{z}^t \rangle - \langle \mathbf{g}^t, \mathbf{z}^* - \mathbf{z}^t \rangle + \frac{\eta}{2\alpha_\Phi} \|\mathbf{g}^t\|_*^2 \right) \\ & = \Psi_1 + \sum_{t=1}^T \left( \frac{1}{2} V(\mathbf{z}^*) - V(\mathbf{z}^t) - \frac{1}{2} \langle \nabla V(\mathbf{z}^t), \mathbf{z}^* - \mathbf{z}^t \rangle + \frac{\eta}{2\alpha_\Phi} \|\mathbf{g}^t\|_*^2 \right) \\ & \leq \frac{1}{\eta} D_\Phi(\mathbf{x}^* \|\mathbf{x}^1) + \frac{\eta}{2\alpha_\Phi} \sum_{t=1}^T \|\mathbf{g}^t\|_*^2. \end{aligned} \quad (13)$$

The third inequality holds because of the (1, 2)-concavity of  $V$ , i.e.,  $V(\mathbf{z}^*) - 2V(\mathbf{z}^t) - \langle \nabla V(\mathbf{z}^t), \mathbf{z}^* - \mathbf{z}^t \rangle \leq 0$ . The theorem follows. □

Combining Theorem 5 and Theorem 2, we obtain the following result.

**Theorem 6** Using Algorithm 1 with the specification that in line 8, replace  $f^t(\mathbf{x}^t)$  by  $\frac{1}{2}v(S^t) - \sum_{j \in S^t} p_j$ , one gets a bandit algorithm with the following guarantee.

$$\sum_{t=1}^T \left( v(S^t) - \sum_{j \in S^t} p_j^t \right) \geq \max_{S \subset [m]} \sum_{t=1}^T \left( \frac{1}{2} v(S) - \sum_{j \in S} p_j^t \right) - O(m^{3/2} (\log T)^{3/2} (\log \log T) \sqrt{T}).$$

## 7 Conclusion

In this paper, we have introduced a framework to design efficient online learning algorithms. Apart of standard regularity requirements (such as compact convex domain, Lipschitz, etc), a new crucial property is the  $(\lambda, \mu)$ -concavity. Designing efficient online learning algorithms is now reduced to constructing  $(\lambda, \mu)$ -concave offline algorithms (also with other standard regularity conditions). We show the applicability of the framework through applications in auction design. Due to the simplicity of the conditions, we hope that our approach would be useful in designing efficient online algorithms with approximate regret bounds for different problems.

## References

- [1] Jacob D. Abernethy, Elad Hazan, and Alexander Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. In *Proc. 21st Annual Conference on Learning Theory (COLT)*, pages 263–274, 2008.
- [2] Naman Agarwal, Alon Gonen, and Elad Hazan. Learning in non-convex games with an optimization oracle. In *Proc. 32nd Conference on Learning Theory*, volume 99, pages 18–29, 2019.
- [3] Baruch Awerbuch and Robert Kleinberg. Online linear optimization and adaptive routing. *Journal of Computer and System Sciences*, 74(1):97–114, 2008.
- [4] Nikhil Bansal and Anupam Gupta. Potential-Function Proofs for First-Order Methods. *arXiv:1712.04581*, 2017.
- [5] Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- [6] Dimitri Bertsekas and Angelia Nedic. Convex analysis and optimization. 2003.
- [7] Avrim Blum and Jason D Hartline. Near-optimal online auctions. In *Proc. 16th Symposium on Discrete Algorithms*, pages 1156–1163, 2005.
- [8] George W Brown. Iterative solution of games by fictitious play. *Activity analysis of production and allocation*, 13(1):374–376, 1951.
- [9] Sébastien Bubeck and Ronen Eldan. The entropic barrier: a simple and optimal universal self-concordant barrier. *Mathematics of Operations Research*, 2018.
- [10] Sébastien Bubeck, Nicolo Cesa-Bianchi, and Sham Kakade. Towards minimax policies for online linear optimization with bandit feedback. In *Annual Conference on Learning Theory*, volume 23, pages 41–1, 2012.
- [11] Sébastien Bubeck, Yin Tat Lee, and Ronen Eldan. Kernel-based methods for bandit convex optimization. In *Proc. 49th Symposium on Theory of Computing*, pages 72–85, 2017.
- [12] 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, 2015.
- [13] Varsha Dani, Sham M Kakade, and Thomas P Hayes. The price of bandit information for online optimization. In *Advances in Neural Information Processing Systems*, pages 345–352, 2008.
- [14] Constantinos Daskalakis and Vasilis Syrgkanis. Learning in auctions: Regret is hard, envy is easy. In *57th Annual Symposium on Foundations of Computer Science*, pages 219–228, 2016.
- [15] Miroslav Dudik, Nika Haghtalab, Haipeng Luo, Robert E Schapire, Vasilis Syrgkanis, and Jennifer Wortman Vaughan. Oracle-efficient online learning and auction design. In *Proc. 58th Symposium on Foundations of Computer Science (FOCS)*, pages 528–539, 2017.
- [16] Shaddin Dughmi, Tim Roughgarden, and Qiqi Yan. From convex optimization to randomized mechanisms: toward optimal combinatorial auctions. In *Proc. 43rd ACM Symposium on Theory of Computing*, pages 149–158, 2011.

- [17] Abraham D Flaxman, Adam Tauman Kalai, and H Brendan McMahan. Online convex optimization in the bandit setting: gradient descent without a gradient. In *Proc. 16th Symposium on Discrete Algorithms*, pages 385–394, 2005.
- [18] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139, 1997.
- [19] Drew Fudenberg and David K Levine. Consistency and cautious fictitious play. *Journal of Economic Dynamics and Control*, 19(5-7):1065–1089, 1995.
- [20] Drew Fudenberg and David K Levine. *The theory of learning in games*. MIT press, 1998.
- [21] James Hannan. Approximation to bayes risk in repeated play. *Contributions to the Theory of Games*, 3:97–139, 1957.
- [22] Hamed Hassani, Mahdi Soltanolkotabi, and Amin Karbasi. Gradient methods for submodular maximization. In *Advances in Neural Information Processing Systems*, pages 5841–5851, 2017.
- [23] Elad Hazan. Introduction to online convex optimization. *Foundations and Trends® in Optimization*, 2(3-4):157–325, 2016.
- [24] Elad Hazan and Satyen Kale. Online submodular minimization. *Journal of Machine Learning Research*, 13:2903–2922, 2012.
- [25] Elad Hazan and Yuanzhi Li. An optimal algorithm for bandit convex optimization. *arXiv preprint arXiv:1603.04350*, 2016.
- [26] Ramesh Johari. Lecture 6: Fictitious play. Lecture notes, April 23 2007.
- [27] Adam Kalai and Santosh Vempala. Efficient algorithms for online decision problems. *Journal of Computer and System Sciences*, 71(3):291–307, 2005.
- [28] Robert Kleinberg and Tom Leighton. The value of knowing a demand curve: Bounds on regret for online posted-price auctions. In *Proc. 44th Symposium on Foundations of Computer Science*, pages 594–605, 2003.
- [29] Robert D Kleinberg. Nearly tight bounds for the continuum-armed bandit problem. In *Advances in Neural Information Processing Systems*, pages 697–704, 2005.
- [30] Nick Littlestone and Manfred K Warmuth. The weighted majority algorithm. *Information and computation*, 108(2):212–261, 1994.
- [31] Thodoris Lykouris, Vasilis Syrgkanis, and Éva Tardos. Learning and efficiency in games with dynamic population. In *Proc. 27th Symposium on Discrete algorithms*, pages 120–129, 2016.
- [32] Roger B Myerson. Optimal auction design. *Mathematics of Operations Research*, 6(1):58–73, 1981.
- [33] Hariharan Narayanan and Alexander Rakhlin. Random walk approach to regret minimization. In *Advances in Neural Information Processing Systems*, pages 1777–1785, 2010.
- [34] Yurii Nesterov and Arkadii Nemirovskii. *Interior-point polynomial algorithms in convex programming*, volume 13. SIAM, 1994.
- [35] Tim Roughgarden. Intrinsic robustness of the price of anarchy. *Journal of the ACM (JACM)*, 62(5):32, 2015.
- [36] Tim Roughgarden. The price of anarchy in games of incomplete information. *ACM Transactions on Economics and Computation*, 3(1):6, 2015.
- [37] Tim Roughgarden and Joshua R Wang. Minimizing regret with multiple reserves. In *Proc. 2016 ACM Conference on Economics and Computation*, pages 601–616, 2016.
- [38] Tim Roughgarden, Vasilis Syrgkanis, and Eva Tardos. The price of anarchy in auctions. *Journal of Artificial Intelligence Research*, 59:59–101, 2017.

- [39] Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194, 2012.
- [40] Vasilis Syrgkanis and Eva Tardos. Composable and efficient mechanisms. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 211–220. ACM, 2013.
- [41] Vasilis Syrgkanis, Akshay Krishnamurthy, and Robert Schapire. Efficient algorithms for adversarial contextual learning. In *International Conference on Machine Learning*, pages 2159–2168, 2016.