
Bandit Convex Optimization: Towards Tight Bounds

Elad Hazan

Technion—Israel Institute of Technology
Haifa 32000, Israel
ehazan@ie.technion.ac.il

Kfir Y. Levy

Technion—Israel Institute of Technology
Haifa 32000, Israel
kfiryl@tx.technion.ac.il

Abstract

Bandit Convex Optimization (BCO) is a fundamental framework for decision making under uncertainty, which generalizes many problems from the realm of online and statistical learning. While the special case of linear cost functions is well understood, a gap on the attainable regret for BCO with *nonlinear* losses remains an important open question. In this paper we take a step towards understanding the best attainable regret bounds for BCO: we give an efficient and near-optimal regret algorithm for BCO with strongly-convex and smooth loss functions. In contrast to previous works on BCO that use time invariant exploration schemes, our method employs an exploration scheme that shrinks with time.

1 Introduction

The power of Online Convex Optimization (OCO) framework is in its ability to generalize many problems from the realm of online and statistical learning, and supply universal tools to solving them. Extensive investigation throughout the last decade has yield efficient algorithms with worst case guarantees. This has lead many practitioners to embrace the OCO framework in modeling and solving real world problems.

One of the greatest challenges in OCO is finding tight bounds to the problem of Bandit Convex Optimization (BCO). In this “bandit” setting the learner observes the loss function only at the point that she has chosen. Hence, the learner has to balance between exploiting the information she has gathered and between exploring the new data. The seminal work of [5] elegantly resolves this “exploration-exploitation” dilemma by devising a combined explore-exploit gradient descent algorithm. They obtain a bound of $\mathcal{O}(T^{3/4})$ on the expected regret for the general case of an adversary playing bounded and Lipschitz-continuous convex losses.

In this paper we investigate the BCO setting assuming that the adversary is limited to inflicting strongly-convex and smooth losses and the player may choose points from a *constrained* decision set. In this setting we devise an efficient algorithm that achieves a regret of $\tilde{O}(\sqrt{T})$. This rate is the best possible up to logarithmic factors as implied by a recent work of [13], obtaining a lower bound of $\Omega(\sqrt{T})$ for the same setting.

During our analysis, we develop a full-information algorithm that takes advantage of the strong-convexity of loss functions and uses a self-concordant barrier as a regularization term. This algorithm enables us to perform “shrinking exploration” which is a key ingredient in our BCO algorithm. Conversely, all previous works on BCO use a time invariant exploration scheme.

Setting	Convex	Linear	Smooth	Str.-Convex	Str.-Convex & Smooth
Full-Info.	$\Theta(\sqrt{T})$			$\Theta(\log T)$	
BCO	$\tilde{O}(T^{3/4})$	$\tilde{O}(\sqrt{T})$	$\tilde{O}(T^{2/3})$		$\tilde{O}(\sqrt{T})$ [Thm. 10]
	$\Omega(\sqrt{T})$				

Table 1: Known regret bounds in the Full-Info./ BCO setting. Our new result is highlighted, and improves upon the previous $\tilde{O}(T^{2/3})$ bound.

This paper is organized as follows. In Section 2 we introduce the BCO setting, and review necessary preliminaries regarding self-concordant barriers. Section 3 discusses schemes to perform single-point gradient estimations, then we define first-order online methods and analyze the performance of such methods receiving noisy gradient estimates. Our main result is described and analyzed in Section 4; Section 5 concludes.

1.1 Prior work

For BCO with general convex loss functions, almost simultaneously to [5], a bound of $\mathcal{O}(T^{3/4})$ was also obtained by [9] for the setting of Lipschitz-continuous convex losses. Conversely, the best known lower bound for this problem is $\Omega(\sqrt{T})$ proved for the easier full-information setting.

In case the adversary is limited to using linear losses, it can be shown that the player does not “pay” for exploration; this property was used by [4] to devise the Geometric Hedge algorithm that achieves an optimal regret rate of $\tilde{O}(\sqrt{T})$. Later [1], inspired by interior point methods, devised the first *efficient* algorithm that attains the same nearly-optimal regret rate for this setup of bandit linear optimization.

For some special classes of nonlinear convex losses, there are several works that lean on ideas from [5] to achieve improved upper bounds for BCO. In the case of convex and smooth losses [11] attained an upper bound of $\tilde{O}(T^{2/3})$. The same regret rate of $\tilde{O}(T^{2/3})$ was achieved by [2] in the case of strongly-convex losses. For the special case of *unconstrained* BCO with strongly-convex and smooth losses, [2] obtained a regret of $\tilde{O}(\sqrt{T})$. A recent paper by Shamir [13], significantly advanced our understanding of BCO by devising a lower bound of $\Omega(\sqrt{T})$ for the setting of strongly-convex and smooth BCO. The latter implies the tightness of our bound.

A comprehensive survey by Bubeck and Cesa-Bianchi [3], provides a review of the bandit optimization literature in both stochastic and online setting.

2 Setting and Background

Notation: During this paper we denote by $\|\cdot\|$ the ℓ_2 norm when referring to vectors, and use the same notation for the spectral norm when referring to matrices. We denote by \mathbb{B}^n and \mathbb{S}^n the n -dimensional euclidean unit ball and unit sphere, and by $v \sim \mathbb{B}^n$ and $u \sim \mathbb{S}^n$ random variables chosen uniformly from these sets. The symbol \mathcal{I} is used for the identity matrix (its dimension will be clear from the context). For a positive definite matrix $A \succ 0$ we denote by $A^{1/2}$ the matrix B such that $B^\top B = A$, and by $A^{-1/2}$ the inverse of B . Finally, we denote $[N] := \{1, \dots, N\}$.

2.1 Bandit Convex Optimization

We consider a repeated game of T rounds between a player and an adversary, at each round $t \in [T]$

1. player chooses a point $x_t \in \mathcal{K}$.

2. adversary independently chooses a loss function $f_t \in \mathcal{F}$.
3. player suffers a loss $f_t(x_t)$ and receives a feedback \mathbb{F}_t .

In the OCO (Online Convex Optimization) framework we assume that the decision set \mathcal{K} is convex and that all functions in \mathcal{F} are convex. Our paper focuses on adversaries limited to choosing functions from the set $\mathcal{F}_{\sigma,\beta}$; the set of all σ -strongly convex and β -smooth functions.

We also limit ourselves to *oblivious* adversaries where the loss sequence $\{f_t\}_{t=1}^T$ is predetermined and is therefore independent of the player's choices. Mind that in this case the best point in hindsight is also independent of the player's choices. We also assume that the loss functions are defined over the entire space \mathbb{R}^n and are strongly-convex and smooth there; yet the player may only choose points from a constrained set \mathcal{K} .

Let us define the regret of \mathcal{A} , and its regret with respect to a comparator $w \in \mathcal{K}$:

$$\text{Regret}_T^{\mathcal{A}} = \sum_{t=1}^T f_t(x_t) - \min_{w^* \in \mathcal{K}} \sum_{t=1}^T f_t(w^*), \quad \text{Regret}_T^{\mathcal{A}}(w) = \sum_{t=1}^T f_t(x_t) - \sum_{t=1}^T f_t(w)$$

A player aims at minimizing his regret, and we are interested in players that ensure an $o(T)$ regret for any loss sequence that the adversary may choose.

The player learns through the feedback \mathbb{F}_t received in response to his actions. In the full informations setting, he receives the loss function f_t itself as a feedback, usually by means of a gradient oracle - i.e. the decision maker has access to the gradient of the loss function at any point in the decision set. Conversely, in the BCO setting the given feedback is $f_t(x_t)$, i.e., the loss function only at the point that he has chosen; and the player aims at minimizing his *expected regret*, $\mathbf{E}[\text{Regret}_T^{\mathcal{A}}]$.

2.2 Strong Convexity and Smoothness

As mentioned in the last subsection we consider an adversary limited to choosing loss functions from the set $\mathcal{F}_{\sigma,\beta}$, the set of σ -strongly convex and β -smooth functions, here we define these properties.

Definition 1. (Strong Convexity) We say that a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is σ -strongly convex over the set \mathcal{K} if for all $x, y \in \mathcal{K}$ it holds that,

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\sigma}{2} \|x - y\|^2 \quad (1)$$

Definition 2. (Smoothness) We say that a convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is β -smooth over the set \mathcal{K} if the following holds:

$$f(y) \leq f(x) + \nabla f(x)^\top (y - x) + \frac{\beta}{2} \|x - y\|^2, \quad \forall x, y \in \mathcal{K} \quad (2)$$

2.3 Self Concordant Barriers

Interior point methods are polynomial time algorithms to solving *constrained* convex optimization programs. The main tool in these methods is a *barrier function* that encodes the constrained set and enables the use of a fast *unconstrained* optimization machinery. More on this subject can be found in [10].

Let $\mathcal{K} \in \mathbb{R}^n$ be a convex set with a non empty interior $\text{int}(\mathcal{K})$

Definition 3. A function $\mathcal{R} : \text{int}(\mathcal{K}) \rightarrow \mathbb{R}$ is called ν -self-concordant if:

1. \mathcal{R} is three times continuously differentiable and convex, and approaches infinity along any sequence of points approaching the boundary of \mathcal{K} .

2. For every $h \in \mathbb{R}^n$ and $x \in \text{int}(\mathcal{K})$ the following holds:

$$|\nabla^3 \mathcal{R}(x)[h, h, h]| \leq 2(\nabla^2 \mathcal{R}(x)[h, h])^{3/2} \quad \text{and} \quad |\nabla \mathcal{R}(x)[h]| \leq \nu^{1/2}(\nabla^2 \mathcal{R}(x)[h, h])^{1/2}$$

here the third order differential is defined as:

$$\nabla^3 \mathcal{R}(x)[h, h, h] := \frac{\partial^3}{\partial t_1 \partial t_2 \partial t_3} \mathcal{R}(x + t_1 h + t_2 h + t_3 h) \Big|_{t_1=t_2=t_3=0}$$

Our algorithm requires a ν -self-concordant barrier over \mathcal{K} , and its regret depends on $\sqrt{\nu}$. It is well known that any convex set in \mathbb{R}^n admits a $\nu = \mathcal{O}(n)$ such barrier (ν might be much smaller), and that most interesting convex sets admit a self-concordant barrier that is efficiently represented.

The Hessian of a self-concordant barrier induces a local norm at every $x \in \text{int}(\mathcal{K})$, we denote this norm by $\|\cdot\|_x$ and its dual by $\|\cdot\|_x^*$ and define $\forall h \in \mathbb{R}^n$:

$$\|h\|_x = \sqrt{h^\top \nabla^2 \mathcal{R}(x) h}, \quad \|h\|_x^* = \sqrt{h^\top (\nabla^2 \mathcal{R}(x))^{-1} h}$$

we assume that $\nabla^2 \mathcal{R}(x)$ always has a full rank.

The following fact is a key ingredient in the sampling scheme of BCO algorithms [1, 11]. Let \mathcal{R} be self-concordant barrier and $x \in \text{int}(\mathcal{K})$ then the *Dikin Ellipsoide*,

$$W_1(x) := \{y \in \mathbb{R}^n : \|y - x\|_x \leq 1\} \quad (3)$$

i.e. the $\|\cdot\|_x$ -unit ball centered around x , is completely contained in \mathcal{K} .

Our regret analysis requires a bound on $\mathcal{R}(y) - \mathcal{R}(x)$; hence, we will find the following lemma useful:

Lemma 4. Let \mathcal{R} be a ν -self-concordant function over \mathcal{K} , then:

$$\mathcal{R}(y) - \mathcal{R}(x) \leq \nu \log \frac{1}{1 - \pi_x(y)}, \quad \forall x, y \in \text{int}(\mathcal{K})$$

where $\pi_x(y) = \inf\{t \geq 0 : x + t^{-1}(y - x) \in \mathcal{K}\}$, $\forall x, y \in \text{int}(\mathcal{K})$

Note that $\pi_x(y)$ is called the Minkowsky function and it is always in $[0, 1]$. Moreover, as y approaches the boundary of \mathcal{K} then $\pi_x(y) \rightarrow 1$.

3 Single Point Gradient Estimation and Noisy First-Order Methods

3.1 Single Point Gradient Estimation

A main component of BCO algorithms is a randomized sampling scheme for constructing gradient estimates. Here, we survey the previous schemes as well as the more general scheme that we use.

Spherical estimators: Flaxman et al. [5] introduced a method that produces single point gradient estimates through spherical sampling. These estimates are then inserted into a full-information procedure that chooses the next decision point for the player. Interestingly, these gradient estimates are unbiased predictions for the gradients of a *smoothed version function* which we next define.

Let $\delta > 0$, and assume $v \sim \mathbb{B}^n$, the smoothed version of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is defined as follows:

$$\hat{f}(x) = \mathbf{E}[f(x + \delta v)] \quad (4)$$

The next lemma of [5] ties between the gradients of \hat{f} and an estimate based on samples of f :

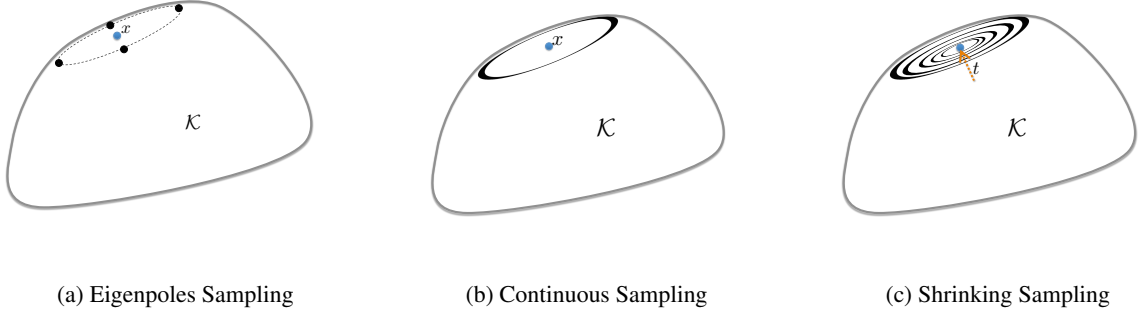


Figure 1: Dikin Ellipsoide Sampling Schemes

Lemma 5. Let $u \sim \mathbb{S}^n$, and consider the smoothed version \hat{f} defined in Equation (4), then the following applies:

$$\nabla \hat{f}(x) = \mathbf{E}\left[\frac{n}{\delta} f(x + \delta u)u\right] \quad (5)$$

Therefore, $\frac{n}{\delta} f(x + \delta u)u$ is an unbiased estimator for the gradients of the smoothed version.

Ellipsoidal estimators: Abernethy et al. [1] introduced the idea of sampling from an ellipsoid (specifically the Dikin ellipsoid) rather than a sphere in the context of BCO. They restricted the sampling to the eigenpoles of the ellipsoid (Fig. 1a). A more general method of sampling continuously from an ellipsoid was introduced in [11] (Fig. 1b). Later we shall introduce a “shrinking-sampling” scheme (Fig. 1c). The following lemma of [11] shows that we can sample f non uniformly over all directions and create an unbiased gradient estimate of a respective smoothed version:

Corollary 6. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuous function, let $A \in \mathbb{R}^{n \times n}$ be invertible, and $v \sim \mathbb{B}^n$, $u \sim \mathbb{S}^n$. Define the smoothed version of f with respect to A :

$$\hat{f}(x) = \mathbf{E}[f(x + Av)] \quad (6)$$

Then the following holds:

$$\nabla \hat{f}(x) = \mathbf{E}[nf(x + Au)A^{-1}u] \quad (7)$$

We prove Corollary 6 in Appendix A.1. Also note that if $A \succ 0$ then $\{Au : u \sim \mathbb{S}^n\}$ is an ellipsoid.

Our next lemma shows that the smoothed version preserves the strong-convexity of f , and that we can measure the distance between \hat{f} and f using the spectral norm of A^2 :

Lemma 7. Consider a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, and a positive definite matrix $A \in \mathbb{R}^{n \times n}$. Let \hat{f} be the smoothed version of f with respect to A as defined in Equation (6). Then the following holds:

- If f is σ -strongly convex then so is \hat{f} .
- If f is convex and β -smooth, and λ_{\max} be the largest eigenvalue of A then:

$$0 \leq \hat{f}(x) - f(x) \leq \frac{\beta}{2} \|A^2\|_2 = \frac{\beta}{2} \lambda_{\max}^2 \quad (8)$$

We prove Lemma 7 in Appendix A.2

Remark: Lemma 7 also holds if we define the smoothed version of f as $\hat{f}(x) = \mathbf{E}_{u \sim \mathbb{S}^n} [f(x + Au)]$ i.e. an average of the original function values over the unit sphere rather than the unit ball as defined in Equation (6). Proof is similar to the one of Lemma 7.

3.2 Noisy First-Order Methods

Our algorithm utilizes a full-information online algorithm, but instead of providing this method with exact gradient values we insert noisy estimates of the gradients. In what follows we define *first-order* online algorithms, and present a lemma that analyses the regret of such algorithm receiving noisy gradients.

Definition 8. (First-Order Online Algorithm) Let \mathcal{A} be an OCO algorithm receiving an arbitrary sequence of differential convex loss functions f_1, \dots, f_T , and providing points $x_1 \leftarrow \mathcal{A}$ and $x_t \leftarrow \mathcal{A}(f_1, \dots, f_{t-1})$. Given that \mathcal{A} requires all loss functions to belong to some set \mathcal{F}_0 . Then \mathcal{A} is called first-order online algorithm if the following holds:

- Adding a linear function to a member of \mathcal{F}_0 remains in \mathcal{F}_0 ; i.e., for every $f \in \mathcal{F}_0$ and $a \in \mathbb{R}^n$ then also $f + a^\top x \in \mathcal{F}_0$
- The algorithm's choices depend only on its gradient values taken in the past choices of \mathcal{A} , i.e. :

$$\mathcal{A}(f_1, \dots, f_{t-1}) = \mathcal{A}(\nabla f_1(x_1), \dots, \nabla f_{t-1}(x_{t-1})), \quad \forall t \in [T]$$

The following is a generalization of Lemma 3.1 from [5]:

Lemma 9. Let w be a fixed point in \mathcal{K} . Let \mathcal{A} be a first-order online algorithm receiving a sequence of differential convex loss functions $f_1, \dots, f_T : \mathcal{K} \rightarrow \mathbb{R}$ (f_{t+1} possibly depending on z_1, \dots, z_t). Where $z_1 \dots z_T$ are defined as follows: $z_1 \leftarrow \mathcal{A}$, $z_t \leftarrow \mathcal{A}(g_1, \dots, g_{t-1})$ where g_t 's are vector valued random variables such that:

$$\mathbf{E}[g_t | z_1, f_1, \dots, z_t, f_t] = \nabla f_t(z_t)$$

Then if \mathcal{A} ensures a regret bound of the form: $\text{Regret}_T^{\mathcal{A}} \leq B_{\mathcal{A}}(\nabla f_1(x_1), \dots, \nabla f_T(x_T))$ in the full information case then, in the case of noisy gradients it ensures the following bound:

$$\mathbf{E}\left[\sum_{t=1}^T f_t(z_t)\right] - \sum_{t=1}^T f_t(w) \leq \mathbf{E}[B_{\mathcal{A}}(g_1, \dots, g_T)]$$

Lemma 9 is proved in Appendix A.3.

4 Main Result

Following is the main theorem of this paper:

Theorem 10. Let \mathcal{K} be a convex set with diameter $\mathcal{D}_{\mathcal{K}}$ and \mathcal{R} be a ν -self-concordant barrier over \mathcal{K} . Then in the BCO setting where the adversary is limited to choosing β -smooth and σ -strongly convex functions and $|f_t(x)| \leq L$, $\forall x \in \mathcal{K}$, then the expected regret of Algorithm 1 with $\eta = \sqrt{\frac{(\nu+2\beta/\sigma) \log T}{2n^2 L^2 T}}$ is upper bounded as

$$\mathbf{E}[\text{Regret}_T] \leq 4nL \sqrt{\left(\nu + \frac{2\beta}{\sigma}\right) T \log T} + 2L + \frac{\beta \mathcal{D}_{\mathcal{K}}^2}{2} = \mathcal{O}\left(\sqrt{\frac{\beta \nu}{\sigma} T \log T}\right)$$

whenever $T/\log T \geq 2(\nu + 2\beta/\sigma)$.

Algorithm 1 BCO Algorithm for Str.-convex & Smooth losses

Input: $\eta > 0, \sigma > 0, \nu$ -self-concordant barrier \mathcal{R}
Choose $x_1 \leftarrow \text{FTARL-}\sigma$
for $t = 1, 2 \dots T$ **do**
 Define $B_t = (\nabla^2 \mathcal{R}(x_t) + \eta \sigma t \mathcal{I})^{-1/2}$
 Draw $u \sim \mathbb{S}^n$
 Play $y_t = x_t + B_t u$
 Observe $f_t(x_t + B_t u)$ and define $g_t = n f_t(x_t + B_t u) B_t^{-1} u$
 Update $x_{t+1} \leftarrow \text{FTARL-}\sigma(g_1, \dots, g_t)$
end for

Note that the sampling scheme defined in Algorithm 1 shrinks the exploration magnitude with time. This “shrinking-exploration” (Fig. 1c) is enabled thanks to the strong-convexity of the loss functions. Our BCO algorithm uses a full-information first-order algorithm denoted FTARL- σ (Follow The Approximate Regularized Leader) which is defined below. This algorithm is a variant of the FTRL methodology as defined in [6] and [12]. The regret bound of the FTARL- σ is stated in the following analysis section.

Algorithm 2 FTARL- σ

Input: $\eta > 0, \nu$ -self concordant barrier \mathcal{R}
Choose $x_1 = \arg \min_{x \in \mathcal{K}} \mathcal{R}(x)$
for $t = 1, 2 \dots T$ **do**
 Receive $\nabla h_t(x_t)$
 Output $x_{t+1} = \arg \min_{x \in \mathcal{K}} \sum_{\tau=1}^t \{ \nabla h_t(x_t)^\top x + \frac{\sigma}{2} \|x - x_\tau\|^2 \} + \frac{1}{\eta} \mathcal{R}(x)$
end for

4.1 Analysis

Here we prove Theorem 10.

Let us decompose the expected regret of Algorithm 1 with respect to $w \in \mathcal{K}$:

$$\begin{aligned} \mathbf{E} [\text{Regret}_T(w)] &:= \sum_{t=1}^T \mathbf{E} [f_t(y_t) - f_t(w)] \\ &= \sum_{t=1}^T \mathbf{E} [f_t(y_t) - f_t(x_t)] \end{aligned} \tag{9}$$

$$+ \sum_{t=1}^T \mathbf{E} [f_t(x_t) - \hat{f}_t(x_t)] \tag{10}$$

$$- \sum_{t=1}^T \mathbf{E} [f_t(w) - \hat{f}_t(w)] \tag{11}$$

$$+ \sum_{t=1}^T \mathbf{E} [\hat{f}_t(x_t) - \hat{f}_t(w)] \tag{12}$$

where expectation is taken with respect to the player’s choices, and \hat{f}_t is defined as

$$\hat{f}_t(x) = \mathbf{E}[f_t(x + B_t v)], \quad \forall x \in \mathcal{K}$$

here $v \sim \mathbb{B}^n$ and the smoothing matrix B_t is defined in Algorithm 1.

The sampling scheme used by Algorithm 1 yields an unbiased gradient estimate g_t of the smoothed version \hat{f}_t , which is then inserted to FTARL- σ (Algorithm 2). We can therefore interpret Algorithm 1 as performing

noisy first-order method (FTARL- σ) over the smoothed versions. The x_t 's in Algorithm 1 are the outputs of FTARL- σ , thus the term in Equation (12) is associated with “exploitation”. The other terms in Equations (9)-(11) measure the cost of sampling away from x_t , and the distance between the smoothed version and the original function, hence these term are associated with “exploration”. In what follows we analyze these terms separately and show that Algorithm 1 achieves $\tilde{O}(\sqrt{T})$ regret.

4.1.1 The Exploration Terms:

We will now show that the following three bounds hold:

$$\mathbf{E}[f_t(y_t) - f_t(x_t)] = \mathbf{E}[\mathbf{E}_u[f_t(x_t + B_t u)] - f_t(x_t)|x_t] \leq \frac{\beta}{2} \mathbf{E}[\|B_t^2\|_2] \leq \frac{\beta}{2} \frac{1}{\eta \sigma t} \quad (13)$$

$$-\mathbf{E}[f_t(w) - \hat{f}_t(w)] = \mathbf{E}[\mathbf{E}[\hat{f}_t(w) - f_t(w)|x_t]] \leq \frac{\beta}{2} \mathbf{E}[\|B_t^2\|_2] \leq \frac{\beta}{2} \frac{1}{\eta \sigma t} \quad (14)$$

$$\mathbf{E}[f_t(x_t) - \hat{f}_t(x_t)] = \mathbf{E}[\mathbf{E}[f_t(x_t) - \hat{f}_t(x_t)|x_t]] \leq 0 \quad (15)$$

First note that $B_t^2 = (\nabla^2 \mathcal{R}(x_t) + \eta \sigma t \mathcal{I})^{-1}$, and since $\nabla^2 \mathcal{R}(x_t) \succ 0$ then $\|B_t^2\|_2 \leq \frac{1}{\eta \sigma t}$.

Next notice that

$$\mathbf{E}[f_t(y_t)|x_t] = \mathbf{E}_{u \sim \mathbb{S}^n}[f_t(x_t + B_t u)|x_t]$$

i.e., $\mathbf{E}[f_t(y_t)|x_t]$ is a smoothed version of f_t around x_t , where the smoothing is over the unit sphere. By the remark following Lemma 7 which extends the lemma to smoothing over the sphere then (13) holds.

Equations (14) and (15) follow directly from Lemma 7.

4.1.2 The Exploitation Term:

Here we show that the following applies:

$$\sum_{t=1}^T \mathbf{E}[\hat{f}_t(x_t) - \hat{f}_t(w)] \leq 2\eta(nL)^2 T + \frac{1}{\eta} \mathcal{R}(w) \quad (16)$$

In order to prove Equation (16) note that $\sum_{t=1}^T \mathbf{E}[\hat{f}_t(x_t) - \hat{f}_t(w)]$ is the expected regret of the FTARL- σ algorithm receiving $\{g_1, \dots, g_t\}$ which are unbiased gradient estimates of the smoothed versions as we will soon show. Since FTARL- σ is a first order-algorithm, Lemma 9 ensures us that the expected regret choosing $x_{t+1} \leftarrow \text{FTARL-}\sigma(g_1, \dots, g_t)$ is $\mathbf{E} B_{\text{FTARL-}\sigma}(g_1, \dots, g_T)$ where $B_{\text{FTARL-}\sigma}(\cdot, \dots, \cdot)$ is the regret bound achieved by FTARL- σ receiving the exact gradients of the smoothed versions.

First let us show that g_t is an unbiased estimators of $\nabla \hat{f}_t(x_t)$, this is immediate by Corollary 6:

$$\mathbf{E}[g_t|x_t] = \mathbf{E}_u[n f_t(x_t + B_t u) B_t^{-1} u|x_t] = \nabla \hat{f}_t(x_t)$$

The next Lemma bounds the regret of FTARL- σ in the full-information setting (exact gradients):

Lemma 11. *Let \mathcal{R} be a self-concordant barrier over a convex set \mathcal{K} , and $\eta > 0$. Consider an online player receiving σ -strongly convex loss functions h_1, \dots, h_T and choosing points according to FTARL- σ (Algorithm 2), and $\eta \|\nabla h_t(x_t)\|_t^* \leq \frac{1}{2}$, $\forall t \in [T]$. Then the player's regret is upper bounded as follows:*

$$\sum_{t=1}^T h_t(x_t) - \sum_{t=1}^T h_t(w) \leq 2\eta \sum_{t=1}^T (\|\nabla h_t(x_t)\|_t^*)^2 + \frac{1}{\eta} \mathcal{R}(w), \quad \forall z \in \mathcal{K}$$

here $(\|a\|_t^*)^2 = a^T (\nabla^2 \mathcal{R}(x_t) + \eta \sigma t \mathcal{I})^{-1} a$

Lemma 11 is proved in Appendix A.4. Note that according to Lemma 7 the smoothed versions \hat{f}_t' s are σ -strongly convex, hence the above lemma applies.

Using the regret bound of the above lemma, and the unbiasedness of the g_t 's, Lemma 9 ensures us:

$$\sum_{t=1}^T \mathbf{E}[\hat{f}_t(x_t) - \hat{f}_t(w)] \leq 2\eta \sum_{t=1}^T \mathbf{E}[(\|g_t\|_t^*)^2] + \frac{1}{\eta} \mathcal{R}(w)$$

Recall that we compete against an oblivious adversary, meaning that $\{f_t\}_{t=1}^T$ are independent of the player's choices then so is $w^* = \arg \min_{w \in \mathcal{K}} \sum_{t=1}^T f_t(w)$, thus Lemma 9 also holds for $\mathbf{E}[\text{Regret}_T(w^*)]$.

Now, taking the t 'th summand of the latter equation and conditioning on x_t we get:

$$\begin{aligned} \mathbf{E}[(\|g_t\|_t^*)^2 | x_t] &= \mathbf{E} \left[n^2 (f_t(x_t + B_t u))^2 u^\top B_t^{-1} (\nabla^2 \mathcal{R}(x_t) + \eta \sigma t \mathcal{I})^{-1} B_t^{-1} u | x_t \right] \\ &\leq (nL)^2 \mathbf{E}_{u \sim \mathbb{S}^n} [u^\top u] = (nL)^2 \end{aligned} \quad (17)$$

which verifies Equation (16). Note that we used $|f_t(x)| \leq L$, and $B_t^{-1} = (\nabla^2 \mathcal{R}(x_t) + \eta \sigma t \mathcal{I})^{1/2}$.

Similarly to Equation (17) we can verify that $\mathbf{E}[\|g_t\|_t^* | x_t] \leq nL$, combined with the choice of η in Algorithm 1 we can verify that the condition $\eta \|g_t\|_t^* \leq 1/2$ required in Lemma 11 is fulfilled for $T/\log T \geq 2(\nu + 2\beta/\sigma)$.

4.1.3 Combining Regret Terms

Combining Equations (9)-(16) we can bound the regret of Algorithm 1 as follows:

$$\mathbf{E}[\text{Regret}_T(w)] \leq \frac{1}{\eta} \frac{\beta}{\sigma} \sum_{t=1}^T \frac{1}{t} + 2\eta(nL)^2 T + \frac{1}{\eta} \mathcal{R}(w) \leq 2\eta(nL)^2 T + \frac{1}{\eta} \left(\mathcal{R}(w) + \frac{2\beta}{\sigma} \log T \right) \quad (18)$$

here we used $\sum_{t=1}^T \frac{1}{t} \leq 1 + \int_{t=1}^T \frac{1}{t} dt \leq 1 + \log T \leq 2 \log T$, $\forall T \geq 3$.

Recall that $x_1 = \arg \min_{x \in \mathcal{K}} \mathcal{R}(x)$ and assume without loss of generality that $\mathcal{R}(x_1) = 0$ (we can always add \mathcal{R} a constant), then for a point $w \in \mathcal{K}$ such that $\pi_{x_1}(w) \leq 1 - \frac{1}{T}$ Lemma 4 ensures us that:

$$\mathcal{R}(w) \leq \nu \log T$$

If, on the other hand $\pi_{x_1}(w) \geq 1 - \frac{1}{T}$ then consider $w' = (1 - \frac{1}{T})w + \frac{1}{T}x_1$, and we can bound:

$$\begin{aligned} \mathbf{E}[\text{Regret}_T(w)] &\leq \mathbf{E}[\text{Regret}_T(w')] + \mathbf{E} \left[\sum_{t=1}^T f_t(w) - f_t(w') \right] \leq \mathbf{E}[\text{Regret}_T(w')] + TG \|w - w'\| \\ &\leq 2\eta(nL)^2 T + \frac{1}{\eta} \left(\nu + \frac{2\beta}{\sigma} \right) \log T + 2L + \frac{\beta \mathcal{D}_{\mathcal{K}}^2}{2} \end{aligned}$$

here we used the fact that β -smooth convex functions bounded by L over a convex set with diameter $\mathcal{D}_{\mathcal{K}}$ are G -Lipschitz with $G \leq \frac{2L}{\mathcal{D}_{\mathcal{K}}} + \frac{\beta \mathcal{D}_{\mathcal{K}}}{2}$, next we used $\|w - w'\| = \frac{1}{T} \|w - x_1\| \leq \frac{\mathcal{D}_{\mathcal{K}}}{T}$ and finally $\mathcal{R}(w') \leq \nu \log T$.

Taking $\eta = \sqrt{\frac{(\nu + 2\beta/\sigma) \log T}{2n^2 L^2 T}}$ completes the proof of Theorem 10.

Correctness: Note that the Dikin ellipsoid around a point $x_t \in \mathcal{K}$ is defined by $\{x_t + (\nabla^2 \mathcal{R}(x_t))^{-1/2} u, u \in \mathbb{B}^n\}$; our algorithm chooses points from the set $\{x_t + (\nabla^2 \mathcal{R}(x_t) + \eta \sigma t \mathcal{I})^{-1/2} u, u \in \mathbb{S}^n\}$ which is inside the Dikin ellipsoid and therefore belong to \mathcal{K} (recall that the Dikin Ellipsoid is always in \mathcal{K}).

5 Summary and open questions

We have presented an efficient algorithm that attains near optimal regret for the setting of BCO with strongly-convex and smooth losses, advancing our understanding of optimal regret rates for bandit learning.

Perhaps the most important question in bandit learning remains the resolution of the attainable regret bounds for smooth but non-strongly-convex, or vice versa, and generally convex cost functions (see Table 1). Ideally, this should be accompanied by an efficient algorithm, although understanding the optimal rates up to polylogarithmic factors would be a significant advancement by itself.

References

- [1] Jacob Abernethy, Elad Hazan, and Alexander Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. In *COLT*, pages 263–274, 2008.
- [2] Alekh Agarwal, Ofer Dekel, and Lin Xiao. Optimal algorithms for online convex optimization with multi-point bandit feedback. In *COLT*, pages 28–40, 2010.
- [3] Sébastien Bubeck and Nicolo Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.
- [4] Varsha Dani, Thomas P. Hayes, and Sham Kakade. The price of bandit information for online optimization. In *NIPS*, 2007.
- [5] Abraham Flaxman, Adam Tauman Kalai, and H. Brendan McMahan. Online convex optimization in the bandit setting: gradient descent without a gradient. In *SODA*, pages 385–394, 2005.
- [6] Elad Hazan. A survey: The convex optimization approach to regret minimization. In Suvrit Sra, Sebastian Nowozin, and Stephen J. Wright, editors, *Optimization for Machine Learning*, pages 287–302. MIT Press, 2011.
- [7] Elad Hazan, Amit Agarwal, and Satyen Kale. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 69(2-3):169–192, 2007.
- [8] Adam Tauman Kalai and Santosh Vempala. Efficient algorithms for online decision problems. *J. Comput. Syst. Sci.*, 71(3):291–307, 2005.
- [9] Robert D Kleinberg. Nearly tight bounds for the continuum-armed bandit problem. In *NIPS*, volume 17, pages 697–704, 2004.
- [10] Arkadii Nemirovskii. Interior point polynomial time methods in convex programming. Lecture Notes, 2004.
- [11] Ankan Saha and Ambuj Tewari. Improved regret guarantees for online smooth convex optimization with bandit feedback. In *AISTATS*, pages 636–642, 2011.
- [12] Shai Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends in Machine Learning*, 4(2):107–194, 2011.
- [13] Ohad Shamir. On the complexity of bandit and derivative-free stochastic convex optimization. In *Conference on Learning Theory*, pages 3–24, 2013.
- [14] Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *ICML*, pages 928–936, 2003.

A Proofs

A.1 Proof of Corollary 6

Let $v \sim \mathbb{B}^n$ and $u \sim \mathbb{S}^n$, and define an auxiliary function $F(x) = f(Ax)$ and its smoothed version

$$\hat{F}(x) = \mathbf{E}_v[F(x + v)]$$

Using the last definitions we can write:

$$\begin{aligned} \mathbf{E}_u[nf(x + Au)A^{-1}u] &= A^{-1}\mathbf{E}_u[nF(A^{-1}x + u)u] \\ &= A^{-1}\nabla\hat{F}(A^{-1}x) = A^{-1}A\nabla\hat{f}(x) \\ &= \nabla\hat{f}(x) \end{aligned}$$

where in the second equality we used Lemma 5 taken from [5]. The third equality follows since:

$$\nabla\hat{F}(x) = A\mathbf{E}_v[\nabla f(Ax + Av)] = A\nabla\hat{f}(Ax)$$

A.2 Proof of Lemma 7

Let us first prove the σ -strong convexity of \hat{f} assuming σ -strong convexity of f . We will use the following equivalent definition of strong convexity:

A function f is σ -strongly convex over a convex set \mathcal{K} if and only if for all $\forall x, y \in \mathcal{K}$ and $\alpha \in [0, 1]$ it holds that,

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y) - \frac{\sigma}{2}\alpha(1 - \alpha)\|x - y\|^2 \quad (19)$$

Now, by definition of \hat{f} new have:

$$\begin{aligned} \hat{f}(\alpha x + (1 - \alpha)y) &= \mathbf{E}_v[f(\alpha x + (1 - \alpha)y + Av)] \\ &= \mathbf{E}_v[f(\alpha(x + Av) + (1 - \alpha)(y + Av))] \\ &\leq \mathbf{E}_v[\alpha f(x + Av) + (1 - \alpha)f(y + Av) - \frac{\sigma}{2}\alpha(1 - \alpha)\|x - y\|^2] \\ &= \alpha\hat{f}(x) + (1 - \alpha)\hat{f}(y) - \frac{\sigma}{2}\alpha(1 - \alpha)\|x - y\|^2, \quad \forall x, y \in \mathcal{K} \end{aligned}$$

in the inequality we used the σ -strong convexity of f .

Now assuming that f is convex and β -smooth we prove the second part of the lemma. The LHS of the Equation (8) follows directly from Jensen's inequality and the convexity of f :

$$\hat{f}(x) = \mathbf{E}_v[f(x + Av)] \geq f(\mathbf{E}_v[x + Av]) = f(x)$$

Since A is positive definite, we can write its SVD decomposition as

$$A = W^\top \Lambda W$$

where W is an orthonormal matrix and Λ is diagonal. Letting $\lambda_{\max} := \max_{i \in [n]} \Lambda_{ii}$ be the largest eigenvalue of A , then the RHS of Equation (8) can be written as follows:

$$\begin{aligned}
\hat{f}(x) - f(x) &= \mathbf{E}_v[f(x + Av) - f(x)] \leq \mathbf{E}_v[\nabla f(x)^\top Av + \frac{\beta}{2} \|Av\|^2] \\
&= \frac{\beta}{2} \mathbf{E}_v[v^\top A^\top Av] \\
&= \frac{\beta}{2} \mathbf{E}_v[v^\top W^\top \Lambda^2 W v] \\
&= \frac{\beta}{2} \mathbf{E}_v[v^\top \Lambda^2 v] \\
&\leq \frac{\beta}{2} \lambda_{\max}^2 \mathbf{E}_v[v^\top v] \leq \frac{\beta}{2} \lambda_{\max}^2
\end{aligned}$$

in the first inequality we used the β smoothness of f , the fourth line follows since for an orthonormal matrix W and $v \sim \text{Uni}(B^n)$ then also $Wv \sim \text{Uni}(B^n)$. The fifth line uses $\|v\|^2 \leq 1, \forall v \in \mathbb{B}^n$.

A.3 Proof of Lemma 9

Define the functions $h_t : \mathcal{K} \rightarrow \mathbb{R}$ as follows:

$$h_t(x) = f_t(x) + \xi_t^\top x, \text{ where } \xi_t = g_t - \nabla f_t(z_t)$$

Note that:

$$\nabla h_t(z_t) = \nabla f_t(z_t) + g_t - \nabla f_t(z_t) = g_t$$

Therefore applying deterministically a first-order method \mathcal{A} on the random functions h_t is equivalent to applying \mathcal{A} on a stochastic-first order approximation of the deterministic functions f_t . Thus by the full-information regret bound of \mathcal{A} we have:

$$\sum_{t=1}^T h_t(z_t) - \sum_{t=1}^T h_t(w) \leq B_{\mathcal{A}}(g_1, \dots, g_T) \quad (20)$$

Also note that:

$$\begin{aligned}
\mathbf{E}[h_t(z_t)] &= \mathbf{E}[f_t(z_t)] + \mathbf{E}[\xi_t^\top z_t] = \mathbf{E}[f_t(z_t)] + \mathbf{E}[\mathbf{E}[\xi_t^\top z_t | z_1, f_1, \dots, z_t, f_t]] \\
&= \mathbf{E}[f_t(z_t)] + \mathbf{E}[\mathbf{E}[\xi_t | z_1, f_1, \dots, z_t, f_t]^\top z_t] = \mathbf{E}[f_t(z_t)]
\end{aligned}$$

where we used $\mathbf{E}[\xi_t | z_1, f_1, \dots, z_t, f_t] = 0$. Similarly, since $w \in \mathcal{K}$ is fixed we have that $\mathbf{E}[h_t(w)] = f_t(w)$. Taking the expectation of equation (20) the lemma follows.

A.4 Proof of Lemma 11

First note the following lemma due to [14] (proof is found in [7]):

Lemma 12. *Let h_1, \dots, h_T be an arbitrary sequence of loss functions, and let $x_1, \dots, x_T \in \mathcal{K}$. Let, $\tilde{h}_1, \dots, \tilde{h}_T$ be a sequence of loss function that satisfy $\tilde{h}_t(x_t) = h_t(x_t)$, and $\tilde{h}_t(x) \leq h_t(x)$ for all $x \in \mathcal{K}$. Then*

$$\sum_{t=1}^T h_t(x_t) - \min_{x \in \mathcal{K}} \sum_{t=1}^T h_t(x) \leq \sum_{t=1}^T \tilde{h}_t(x_t) - \min_{x \in \mathcal{K}} \sum_{t=1}^T \tilde{h}_t(x).$$

Since we assume the function h_t to be σ -strongly convex then for a point $x_t \in \mathcal{K}$ the following function fulfills the conditions of Lemma 12

$$\tilde{h}_t(x) = h_t(x_t) + \nabla_t^\top (x - x_t) + \frac{\sigma}{2} \|x - x_t\|^2$$

where we used the notation $\nabla_t := \nabla h_t(x_t)$. Thus, it is sufficient to prove the upper bound on the regret for the approximate loss functions $\{\tilde{h}_t\}_{t=1}^T$.

Note that we can write FTARL- σ (Algorithm 2) as follows:

$$\begin{aligned} x_{t+1} &= \arg \min_{x \in \mathcal{K}} \sum_{\tau=1}^t \tilde{h}_\tau(x) + \frac{1}{\eta} \mathcal{R}(x) = \arg \min_{x \in \mathcal{K}} \eta \sum_{\tau=1}^t \left\{ \nabla_\tau^\top x + \frac{\sigma}{2} \|x - x_\tau\|^2 \right\} + \mathcal{R}(x) \\ &= \arg \min_{x \in \mathcal{K}} \eta \sum_{\tau=1}^t \nabla_\tau^\top x + \left(\mathcal{R}(x) + \frac{\eta}{2} \sigma \sum_{\tau=1}^t \|x - x_\tau\|^2 \right) \\ &= \arg \min_{x \in \mathcal{K}} \eta \sum_{\tau=1}^t \nabla_\tau^\top x + \mathcal{R}_t(x) \end{aligned}$$

where $\mathcal{R}_t(x) = \mathcal{R}(x) + \frac{\eta}{2} \sigma \sum_{\tau=1}^t \|x - x_\tau\|^2$. Thus, the FTARL- σ algorithm finds the best point that minimizes the past sum of losses and a regularization term.

Bounding the differences $\nabla_t^\top (x_t - x_{t+1})$ is useful in bounding the regret of FTARL- σ , as seen in the next lemma due to [8] (proof can be found in [6] or in [12]):

Lemma 13. *Let a regularizer function \mathcal{R} , and h_1, \dots, h_T be a sequence of convex cost functions and let $x_t = \arg \min_{x \in \mathcal{K}} \sum_{\tau=1}^{t-1} h_\tau(x) + \eta^{-1} \mathcal{R}(x)$, Then:*

$$\sum_{t=1}^T h_t(x_t) - \sum_{t=1}^T h_t(v) \leq \sum_{t=1}^T \nabla_t^\top (x_t - x_{t+1}) + \eta^{-1} (\mathcal{R}(v) - \mathcal{R}(x_1)), \quad \forall v \in \mathcal{K}$$

In what follows we will bound the differences $\nabla_t^\top (x_t - x_{t+1})$ and use Lemma 13 the regret of FTARL- σ : Given x_t define the following norm and its dual with respect to \mathcal{R}_t :

$$\begin{aligned} \|z\|_t &= \sqrt{z^\top \nabla^2 \mathcal{R}_t(x_t) z} \\ \|z\|_t^* &= \sqrt{z^\top \nabla^2 \mathcal{R}_t(x_t)^{-1} z} \end{aligned}$$

By Holder's inequality we can bound:

$$\nabla_t^\top (x_t - x_{t+1}) \leq \|\nabla_t\|_t^* \|x_t - x_{t+1}\|_t$$

Define the objective of FTARL- σ as: $\Phi_t(x) = \eta \sum_{\tau=1}^t \nabla_\tau^\top x + \mathcal{R}_t(x)$, and define the newton decrement of Φ_t at x as:

$$\lambda(x, \Phi_t) = \|\nabla \Phi_t(x)\|_x^* = \|\nabla^2 \Phi_t(x)^{-1} \nabla \Phi_t(x)\|_x$$

The following lemma helps us measuring the distance between x and the global minimizer of Φ_t :

Lemma 14. *Let g be self-concordant, than if $\lambda(x, g) \leq 1/2$ the following applies:*

$$\|x - \arg \min_x g\|_x \leq 2\lambda(x, g)$$

where $\|z\|_x = \sqrt{z^\top \nabla^2 g(x) z}$.

Using the above lemma together with $x_{t+1} := \arg \min \Phi_t$ and $\nabla^2 \Phi_t = \nabla^2 \mathcal{R}_t$, we get:

$$\|x_t - x_{t+1}\|_t \leq 2\lambda(x, \Phi_t) \leq 2\eta \|\nabla_t\|_t^*$$

We used $\nabla \Phi_t(x_t) = \nabla_t$, which follows since $\Phi_t(x) = \Phi_{t-1}(x) + \eta \nabla_t^\top x + \frac{\eta}{2} \sigma \|x - x_t\|^2$. Using Lemma 13 and assuming that $\eta \|\nabla_t\|_t^* \leq 1/2$ for every $t \in 1, \dots, T$, we have:

$$\text{Regret}_{\text{FTARL}-\sigma}(u) \leq 2\eta \sum_{t=1}^T (\|\nabla_t\|_t^*)^2 + \frac{1}{\eta} \mathcal{R}_0(u)$$

Recalling $(\|z\|_t^*)^2 = z^\top \nabla^2 \mathcal{R}_t(x_t)^{-1} z = z^\top (\nabla^2 \mathcal{R}(x_t) + \eta \sigma t \mathcal{I})^{-1} z$, establishes Lemma 11.