

# Optimistic Bandit Convex Optimization

**Authored by:**

Mehryar Mohri  
Scott Yang

## **Abstract**

We introduce the general and powerful scheme of predicting information re-use in optimization algorithms. This allows us to devise a computationally efficient algorithm for bandit convex optimization with new state-of-the-art guarantees for both Lipschitz loss functions and loss functions with Lipschitz gradients. This is the first algorithm admitting both a polynomial time complexity and a regret that is polynomial in the dimension of the action space that improves upon the original regret bound for Lipschitz loss functions, achieving a regret of  $\widetilde{O}(T^{11/16}d^{3/8})$ . Our algorithm further improves upon the best existing polynomial-in-dimension bound (both computationally and in terms of regret) for loss functions with Lipschitz gradients, achieving a regret of  $\widetilde{O}(T^{8/13}d^{5/3})$ .

## **1 Paper Body**

Bandit convex optimization (BCO) is a key framework for modeling learning problems with sequential data under partial feedback. In the BCO scenario, at each round, the learner selects a point (or action) in a bounded convex set and observes the value at that point of a convex loss function determined by an adversary. The feedback received is limited to that information: no gradient or any other higher order information about the function is provided to the learner. The learner's objective is to minimize his regret, that is the difference between his cumulative loss over a finite number of rounds and that of the loss of the best fixed action in hindsight. The limited feedback makes the BCO setup relevant to a number of applications, including online advertising. On the other hand, it also makes the problem notoriously difficult and requires the learner to find a careful trade-off between exploration and exploitation. While it has been the subject of extensive study in recent years, the fundamental BCO problem remains one of the most challenging scenarios in machine learning where several questions concerning optimality guarantees remain open.  $e^{5/6}$  is achievable for bounded The original work of Flaxman et al. [2005] showed that a regret of  $O(T e^{3/4})$  for Lipschitz loss functions (the latter bound is also given in

[Kleinberg, loss functions and of  $O(T^{2/3})$ ], both of which are still the best known results given by explicit algorithms. Agarwal et al. [2010] introduced an algorithm that maintains a regret of  $O(T^{2/3})$  for loss functions that are both Lipschitz and strongly convex, which is also still state-of-the-art. For functions that are Lipschitz and also admit Lipschitz gradients, Saha and Tewari [2011] designed an algorithm with a regret of  $O(T^{2/3})$  regret, a result that was recently improved to  $O(T^{5/8})$  by Dekel et al. [2015].  $O(T^{5/8})$

Here, we further improve upon these bounds both in the Lipschitz and Lipschitz gradient settings. By incorporating the novel and powerful idea of predicting information re-use, we introduce an algorithm with a regret bound of  $O(T^{11/16})$  for Lipschitz loss functions. Similarly, our algorithm also achieves with a regret bound of  $O(T^{11/16})$  the best regret guarantee among computationally tractable algorithms for loss functions with Lipschitz gradients. 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

Both bounds admit a relatively mild dependency on the dimension of the action space:  $O(d^{1/3})$ .

We note that the recent remarkable work by [Bubeck et al., 2015, Bubeck and Eldan, 2015] has a regret bound of  $O(T^{1/2})$ , which matches the existence of algorithms that can attain a regret of  $O(T^{1/2})$  lower bound given by Dani et al.. Thus, the dependency of our bounds with respect to  $T$  is not optimal. Furthermore, two recent unpublished manuscripts, [Hazan and Li, 2016] and [Bubeck et al., 2016]. These results, once verified, would be et al., [2016], present algorithms achieving regret  $O(T^{1/2})$  ground-breaking contributions to the literature. However, unlike our algorithms, the regret bound for both of these algorithms admits a large dependency on the dimension  $d$  of the action space: exponential for [Hazan and Li, 2016],  $d^{O(9.5)}$  for [Bubeck et al., 2016]. One hope is that the novel ideas introduced by Hazan and Li [2016] (the application of the ellipsoid method with a restart button and lower convex envelopes) or those by Bubeck et al. [2016] (which also make use of the restart idea but introduces a very original kernel method) could be combined with those presented in this paper to derive algorithms with the most favorable guarantees with respect to both  $T$  and  $d$ . We begin by formally introducing our notation and setup. We then highlight some of the essential ideas in previous work before introducing our new key insight. Next, we give a detailed description of our algorithm for which we prove theoretical guarantees in several settings.

## 2.1

### Preliminaries BCO scenario

The scenario of bandit convex optimization, which dates back to [Flaxman et al., 2005], is a sequential prediction problem on a convex compact domain  $K \subset \mathbb{R}^d$ . At each round  $t \in [1, T]$ , the learner selects a (possibly) randomized action  $x_t \in K$  and incurs the loss  $f_t(x_t)$  based on a convex function  $f_t : K \rightarrow \mathbb{R}$  chosen by the adversary. We assume that the adversary is oblivious, so that the loss functions are independent of the player's actions. The objective of the learner is to minimize his regret with respect to the optimal static action in hindsight, that is, if we denote by  $A$  the learner's randomized algorithm, the following quantity:  $\text{Reg}_T(A) = \mathbb{E} \sum_{t=1}^T f_t(x_t) - \min_{x \in K} \sum_{t=1}^T f_t(x)$ .

$$\begin{matrix} t=1 \\ t=1 \end{matrix}$$

We will denote by  $D$  the diameter of the action space  $K$  in the Euclidean norm:  $D = \sup_{x,y \in K} \|x - y\|_2$ . Throughout this paper, we will often use different induced norms. We will denote by  $\|\cdot\|_A$  the norm induced by a symmetric positive definite (SPD) matrix  $A \succ 0$ , defined for all  $x \in \mathbb{R}^d$  by  $\|x\|_A = \sqrt{x^\top A x}$ . Moreover, we will denote by  $\|\cdot\|_{A^*}$  its dual norm, given by  $\|x\|_{A^*} = \|x\|_A^{-1}$ . To simplify the notation, we will write  $\|x\|$  instead of  $\|x\|_{R(x)}$ , when the convex and twice differentiable function  $R : \text{int}(K) \rightarrow \mathbb{R}$  is clear from the context. Here,  $\text{int}(K)$  is the set interior of  $K$ .

We will consider different levels of regularity for the functions  $f_t$  selected by the adversary. We will always assume that they are uniformly bounded by some constant  $C \geq 0$ , that is  $-f_t(x) \leq C$  for all  $t \in [1, T]$  and  $x \in K$ , and, by shifting the loss functions upwards by at most  $C$ , we will also assume, without loss of generality, that they are non-negative:  $f_t \geq 0$ , for all  $t \in [1, T]$ . Moreover, we will always assume that  $f_t$  is Lipschitz on  $K$  (henceforth denoted  $C_{0,1}(K)$ ):  $\forall t \in [1, T], \forall x, y \in K, |f_t(x) - f_t(y)| \leq L \|x - y\|_2$ . In some instances, we will further assume that the functions admit  $H$ -Lipschitz gradients on the interior of the domain (henceforth denoted  $C_{1,1}(\text{int}(K))$ ):  $\forall t \in [1, T], \forall x, y \in \text{int}(K), \|\nabla f_t(x) - \nabla f_t(y)\|_2 \leq H \|x - y\|_2$ . Since  $f_t$  is convex, it admits a subgradient at any point in  $K$ . We denote by  $g_t$  one element of the subgradient at the point  $x_t \in K$  selected by the learner at round  $t$ . When the losses are  $C_{1,1}$ , the only element of the subgradient is the gradient, and  $g_t = \nabla f_t(x_t)$ . We will use the shorthand  $\sum_{t=1}^T v_t$  to denote the sum of  $T$  vectors  $v_1, \dots, v_T$ . In particular,  $G_T$  will denote the sum of the subgradients  $g_s$  for  $s \in [1, T]$ . Lastly, we will denote by  $B_1(0) = \{x \in \mathbb{R}^d : \|x\|_2 \leq 1\}$  the  $d$ -dimensional Euclidean ball of radius one and by  $\mathbb{S}^{d-1}$  the unit sphere.

## 2.2

### Follow-the-regularized-leader template

A standard algorithm in online learning, both for the bandit and full-information setting is the follow-the-regularized-leader (FTRL) algorithm. At each round, the algorithm selects the action that minimizes the cumulative linearized loss augmented with a regularization term  $R : K \rightarrow \mathbb{R}$ . Thus, the action  $x_{t+1}$  is defined as follows:  $x_{t+1} = \arg\min_{x \in K} \sum_{s=1}^t \langle g_s, x \rangle + R(x)$

where  $\eta \geq 0$  is a learning rate that determines the tradeoff between greedy optimization and regularization. If we had access to the subgradients at each round, then, FTRL with  $R(x) = \frac{\eta}{2} \|x\|_2^2$  and  $\eta = \frac{1}{\sqrt{T}}$  would yield a regret of  $O(\sqrt{dT})$ , which is known to be optimal. But, since we only have access to the loss function values  $f_t(x_t)$  and since the loss functions change at each round, a more refined strategy is needed.

#### 2.2.1 One-point gradient estimates and surrogate losses

One key insight into the bandit convex optimization problem, due to Flaxman et al. [2005], is that the subgradient of a smoothed version of the loss function can be estimated by sampling and rescaling around the point the algorithm originally intended to play. Lemma 1 ([Flaxman et al., 2005, Saha and Tewari, 2011]). Let  $f : K \rightarrow \mathbb{R}$  be an arbitrary function (not necessarily differen-

tiable) and let  $U(\mathbb{B}_1(0))$  denote the uniform distribution over the unit sphere. Then, for any  $\epsilon > 0$  and any SPD matrix  $A \succ 0$ , the function  $f_b$  defined for all  $x \in K$  by  $f_b(x) = \mathbb{E}_{u \sim U(\mathbb{B}_1(0))} [f(x + Au)]$  is differentiable over  $\text{int}(K)$  and, for any  $x \in \text{int}(K)$ ,  $g_b = d f_b(x + Au)A^{-1}u$  is an unbiased estimate of  $\text{rfb}(x)$ :  $\mathbb{E} d f_b(x + Au)A^{-1}u = \text{rfb}(x)$ .  $u \sim U(\mathbb{B}_1(0))$

The result shows that if at each round  $t$  we sample  $u_t \sim U(\mathbb{B}_1(0))$ , define an SPD matrix  $A_t$  and play the point  $y_t = x_t + A_t u_t$  (assuming that  $y_t \in K$ ), then  $g_{b,t} = d f_b(x_t + A_t u_t)A_t^{-1}u_t$  is an unbiased estimate of the gradient of  $f_b$  at the point  $x_t$  originally intended:  $\mathbb{E}[g_{b,t}] = \text{rfb}(x_t)$ . Thus, we can use FTRL with these smoothed gradient estimates:  $x_{t+1} = \arg\min_{x \in K} \sum_{s=1}^t g_{b,s} \cdot x + R(x)$ , at the cost of the approximation error from  $f_t$  to  $f_b$ . Furthermore, the norm of these estimate gradients can be bounded. Lemma 2. Let  $\epsilon > 0$ ,  $u_t \sim U(\mathbb{B}_1(0))$  and  $A_t \succ 0$ , then the norm of  $g_{b,t} = d f_b(x_t + A_t u_t)A_t^{-1}u_t$  can be bounded as follows:  $\|g_{b,t}\| \leq C \sqrt{d} \sqrt{\lambda_{\min}(A_t)}$ .

Proof. Since  $f_t$  is bounded by  $C$ , we can write  $\|g_{b,t}\| \leq C \sqrt{d} \sqrt{\lambda_{\min}(A_t)}$ .

$d \geq 2$

$C \geq 2$   $u_t \sim U(\mathbb{B}_1(0))$   $A_t \succ 0$   $A_t^{-1}u_t$

$d \geq 2$

$C \geq 2$ .

This gives us a bound on the Lipschitz constant of  $f_{b,t}$  in terms of  $d$ ,  $\epsilon$ , and  $C$ . 2.2.2

Self-concordant barrier as regularization

When sampling to derive a gradient estimate, we need to ensure that the point sampled lies within the feasible set  $K$ . A second key idea in the BCO problem, due to Abernethy et al. [2008], is to design ellipsoids that are always contained in the feasible sets. This is done by using tools from the theory of interior-point methods in convex optimization. Definition 1 (Definition 2.3.1 [Nesterov and Nemirovskii, 1994]). Let  $K \subset \mathbb{R}^d$  be closed convex, and let  $\epsilon > 0$ . A  $C \geq 3$  function  $R : \text{int}(K) \rightarrow \mathbb{R}$  is a  $\epsilon$ -self-concordant barrier for  $K$  if for any sequence  $d(z_s)_{s=1}^\infty$  with  $z_s \in \text{int}(K)$ , we have  $R(z_s) \leq 1$ , and if for all  $x \in \text{int}(K)$ , and  $y \in \mathbb{R}^d$ , the following inequalities hold:  $-\epsilon^3 R(x)[y, y, y] \leq 2\|y\|_x^3$ ,  $-\epsilon R(x)[y] \leq \frac{1}{2}\|y\|_x$ .

Since self-concordant barriers are preserved under translation, we will always assume for convenience that  $\min_{x \in K} R(x) = 0$ . Nesterov and Nemirovskii [1994] show that any  $d$ -dimensional closed convex set admits an  $O(d)$ -self-concordant barrier. This allows us to always choose a self-concordant barrier as regularization. We will use several other key properties of self-concordant barriers in this work, all of which are stated precisely in Appendix 7.1.

3

Previous work

The original paper by Flaxman et al. [2005] sampled indiscriminately around spheres and projected  $e^{-T/4}$  for  $C_0,1$  loss functions. back onto the feasible set at each round. This yielded a regret of  $O(\sqrt{T})$ . The follow-up work of Saha and Tewari [2011] showed that for  $C$  loss functions, one can run FTRL with a self-concordant barrier as regularization and sample around the Dikin ellipsoid to attain an  $e^{-T/3}$  improved regret bound of  $O(\sqrt{T})$ .

More recently, Dekel et al. [2015] showed that by averaging the smoothed gradient estimates  $\bar{g}_t = \frac{1}{k+1} \sum_{i=t-k}^t g_i$  and still using the self-concordant barrier as regularization, one can achieve a regret of  $O(\sqrt{P} \sqrt{k+1})$ . Specifically, denote by  $\bar{g}_t = \frac{1}{k+1} \sum_{i=t-k}^t g_i$  the average of the past  $k+1$  incurred gradients, where  $i=0$  if  $t \leq k$ . Then we can play FTRL on these averaged smoothed gradient estimates:  $x_{t+1} = \arg\min_{x \in K} \sum_{i=1}^t \bar{g}_i^\top x + R(x)$ , to attain the better guarantee.

Abernethy and Rakhlin [2009] derive a generic estimate for FTRL algorithms with self-concordant barriers as regularization: Lemma 3 ([Abernethy and Rakhlin, 2009]-Theorem 2.2-2.3). Let  $K$  be a closed convex set in  $\mathbb{R}^d$  and let  $R$  be a  $\gamma$ -self-concordant barrier for  $K$ . Let  $\{g_t\}_{t=1}^T \subset \mathbb{R}^d$  and  $\eta \geq 0$  be such that  $\eta \sum_{t=1}^T \|g_t\|_{K, R}^2 \leq 1/4$  for all  $t \in [1, T]$ . Then, the FTRL update  $x_{t+1} = \arg\min_{x \in K} \sum_{i=1}^t g_i^\top x + R(x)$  admits the following guarantees:  $\sum_{t=1}^T \ell_t(x_{t+1}) - \min_{x \in K} \sum_{t=1}^T \ell_t(x) \leq \frac{1}{\eta} \sum_{t=1}^T \|g_t\|_{K, R}^2 + R(x^*)$ .

By Lemma 2, if we use FTRL with smoothed gradients, then the upper bound in this lemma can be further bounded by  $\sum_{t=1}^T \ell_t(x_{t+1}) - \min_{x \in K} \sum_{t=1}^T \ell_t(x) \leq \frac{1}{\eta} \sum_{t=1}^T \|g_t\|_{K, R}^2 + R(x^*)$ . Furthermore, the regret is then bounded by the sum of this upper bound and the cost of approximating  $f_t$  with  $\tilde{f}_t$ . On the other hand, Dekel et al. [2015] showed that if we used FTRL with averaged smoothed gradients instead, then the upper bound in this lemma can be bounded as  $\sum_{t=1}^T \ell_t(x_{t+1}) - \min_{x \in K} \sum_{t=1}^T \ell_t(x) \leq \frac{1}{\eta} \sum_{t=1}^T \|g_t\|_{K, R}^2 + 2D \sqrt{L} \sqrt{k+1} \sum_{t=1}^T \|g_t\|_{K, R} + R(x^*)$ . The extra factor  $(k+1)$  in the denominator, at the cost of now approximating  $f_t$  with  $\tilde{f}_t$ , is what contributes to their improved regret result.

In general, finding surrogate losses that can both be approximated accurately and admit only a mild variance is a delicate task, and it is not clear how the constructions presented above can be improved.

#### 4.4.1

##### Algorithm Predicting the predictable

Rather than designing a newer and better surrogate loss, our strategy will be to exploit the structure of the current state-of-the-art method. Specifically, we draw upon the technique of predictable sequences from [Rakhlin and Sridharan, 2013]. The idea here is to allow the learner to preemptively “guess” the gradient at the next step and optimize for this in the FTRL update. Specifically, if  $\tilde{g}_{t+1}$  is an estimate of the time  $t+1$  gradient  $g_{t+1}$  based on information up to time  $t$ , then the learner should play:  $x_{t+1} = \arg\min_{x \in K} \sum_{i=1}^t g_i^\top x + \tilde{g}_{t+1}^\top x + R(x)$ .

This optimistic FTRL algorithm admits the following guarantee: Lemma 4 (Lemma 1 [Rakhlin and Sridharan, 2013]). Let  $K$  be a closed convex set in  $\mathbb{R}^d$ , and let  $R$  be a  $\gamma$ -self-concordant barrier for  $K$ . Let  $\{g_t\}_{t=1}^T \subset \mathbb{R}^d$  and  $\eta \geq 0$  such that  $\eta \sum_{t=1}^T \|g_t\|_{K, R}^2 \leq 1/4$  for all  $t \in [1, T]$ . Then the FTRL update  $x_{t+1} = \arg\min_{x \in K} \sum_{i=1}^t g_i^\top x + \tilde{g}_{t+1}^\top x + R(x)$  admits the following guarantee:  $\sum_{t=1}^T \ell_t(x_{t+1}) - \min_{x \in K} \sum_{t=1}^T \ell_t(x) \leq \frac{1}{\eta} \sum_{t=1}^T \|g_t\|_{K, R}^2 + R(x^*) + 2D \sqrt{L} \sqrt{k+1} \sum_{t=1}^T \|g_t\|_{K, R} + R(x^*)$ .

This optimistic FTRL algorithm admits the following guarantee: Lemma 4 (Lemma 1 [Rakhlin and Sridharan, 2013]). Let  $K$  be a closed convex set in  $\mathbb{R}^d$ , and let  $R$  be a  $\gamma$ -self-concordant barrier for  $K$ . Let  $\{g_t\}_{t=1}^T \subset \mathbb{R}^d$  and  $\eta \geq 0$  such that  $\eta \sum_{t=1}^T \|g_t\|_{K, R}^2 \leq 1/4$  for all  $t \in [1, T]$ . Then the FTRL update  $x_{t+1} = \arg\min_{x \in K} \sum_{i=1}^t g_i^\top x + \tilde{g}_{t+1}^\top x + R(x)$  admits the following guarantee:  $\sum_{t=1}^T \ell_t(x_{t+1}) - \min_{x \in K} \sum_{t=1}^T \ell_t(x) \leq \frac{1}{\eta} \sum_{t=1}^T \|g_t\|_{K, R}^2 + R(x^*) + 2D \sqrt{L} \sqrt{k+1} \sum_{t=1}^T \|g_t\|_{K, R} + R(x^*)$ .

$\arg\min_{\mathbf{x}} \sum_{t=1}^T \ell_t(\mathbf{x}) + R(\mathbf{x})$  admits the following guarantee:  $\sum_{t=1}^T \ell_t(\mathbf{x}_t) - \min_{\mathbf{x}} \sum_{t=1}^T \ell_t(\mathbf{x}) \leq \frac{1}{k+1} \sum_{t=1}^T \|\mathbf{g}_t\|^2 + R(\mathbf{x})$ . In general, it is not clear what would be a good prediction candidate. Indeed, this is why Rakhlin and Sridharan [2013] called this algorithm an ‘‘optimistic’’ FTRL. However, notice that if we elect to play the averaged smoothed losses as in [Dekel et al., 2015], then the update at each time is  $\mathbf{g}_{t+1} = \frac{1}{k+1} \sum_{i=0}^k \mathbf{g}_{t-i}$ . This implies that the time  $t+1$  gradient is  $\mathbf{g}_{t+1} = \frac{1}{k+1} \sum_{i=0}^k \mathbf{g}_{t-i}$ , which includes the smoothed gradients from time  $t+1$  down to time  $t-k$ . The key insight here is that at time  $t$ , all but the  $(t-k)$ -th gradient are known!

This means that if we predict  $\mathbf{g}_{t+1}$

$$\mathbf{g}_{t+1} =$$

$$\frac{1}{k+1} \sum_{i=0}^k \mathbf{g}_{t-i}$$

$$\mathbf{g}_t$$

$$\frac{1}{k+1} \sum_{i=0}^k \mathbf{g}_{t-i} = \mathbf{g}_{t+1}, \quad \mathbf{g}_{t+1} = \frac{1}{k+1} \sum_{i=0}^k \mathbf{g}_{t-i}$$

$$\mathbf{g}_t$$

then the first term in the bound of Lemma 4 will be in terms of  $\mathbf{g}_t$

$$\mathbf{g}_t$$

$$\mathbf{g}_t =$$

$$\frac{1}{k+1} \sum_{i=0}^k \mathbf{g}_{t-i}$$

$$\mathbf{g}_t$$

$$\frac{1}{k+1} \sum_{i=0}^k \mathbf{g}_{t-i}$$

$$\mathbf{g}_t$$

$$\mathbf{g}_t$$

$$=$$

$$\frac{1}{k+1} \sum_{i=0}^k \mathbf{g}_{t-i}$$

In other words, all but the time  $t$  smoothed gradient will cancel out. Essentially, we are predicting the predictable portion of the averaged gradient and guaranteeing that the optimism will pay off. Moreover, where we gained a factor of  $k+1$  in the averaged loss case, we should expect to gain a factor of  $(k+1)^2$  by using this optimistic prediction. Note that this technique of optimistically predicting the variance reduction is widely applicable. As alluded to with the reference to [Schmidt et al., 2013], many variance reduction-type techniques, particularly in stochastic optimization, use historical information in their estimates (e.g. SVRG [Johnson and Zhang, 2013], SAGA [Defazio et al., 2014]). In these cases, it is possible to ‘‘predict’’ the information re-use and improve the convergence rates of each algorithm. 4.2

Description and pseudocode

Here, we give a detailed description of our algorithm, OPTIMISTIC BCO. At each round  $t$ , the algorithm uses a sample  $\mathbf{u}_t$  from the uniform distribution over the unit sphere to define an unbiased estimate of the gradient of  $f_t$ , a smoothed version of the loss function  $f_t$ , as described in Section 2.2.1:  $\mathbf{d}_t = \frac{1}{\sqrt{2}} \nabla f_t(\mathbf{u}_t)$ . Next, the trailing average of these unbiased estimates over a fixed window of length  $k+1$  is computed:  $\mathbf{g}_t = \frac{1}{k+1} \sum_{i=0}^k \mathbf{d}_{t-i}$ . The remaining steps executed at each round coincide with the Follow-the-Regularized-Leader update with a self-concordant barrier used as a regularizer, augmented with an optimistic prediction of the next round’s trailing average.

As described in Section 4.1, all but one of the terms in the trailing average are known and we predict their occurrence:  $k$

$$\begin{aligned} \text{get}+1 &= \\ &1 \times \text{gbt}+1 \text{ i}, k+1 \text{ i}=1 \\ &\downarrow \\ \text{xt}+1 &= \arg\min_{x \in K} \left( \sum_{t=1}^{\text{get}+1} \ell_t(x) + R(x) \right). \end{aligned}$$

Note that Theorem 3 implies that the actual point we play,  $y_t$ , is always a feasible point in  $K$ . Figure 1 presents the pseudocode of the algorithm.

**OPTIMISTIC BCO**( $R, \ell, k, x_1$ )  
1 for  $t = 1$  to  $T$  do  
2  $u_t \leftarrow \text{SAMPLE}(U(\text{B1}(0)))$   
3  $y_t \leftarrow \text{xt} + (r_2 R(x_t))$   
4  $P \leftarrow \text{LAY}(y_t)$   
5  $ft(y_t) \leftarrow R \leftarrow \text{ECEIVE LOSS}(y_t)$   
6  $\text{gbt} \leftarrow \text{ft}(y_t) + (r_2 R(x_t))$   
7  $g_t \leftarrow \text{bt} \text{ i}=0$   
8  $\text{get}+1 \leftarrow \text{bt}+1$   
9  $\text{xt}+1 \leftarrow \arg\min_{x \in K} \left( \sum_{t=1}^{\text{get}+1} \ell_t(x) + R(x) \right)$   
10 return  $\text{ft}(y_t)$

Figure 1: Pseudocode of OPTIMISTIC BCO, with  $R : \text{int}(K) \rightarrow \mathbb{R}$ ,  $x_1 \in K$ .

5  
2  $(0, 1], \ell \in [0, k] \subset \mathbb{Z}$ , and  
Regret guarantees

In this section, we state our main results, which are regret guarantees for OPTIMISTIC BCO in the  $C_{0,1}$  and  $C_{1,1}$  cases. We also highlight the analysis and proofs for each regime.

Main results

The following is our main result for the  $C_{0,1}$  case. Theorem 1 ( $C_{0,1}$  Regret). Let  $K \subset \mathbb{R}^d$  be a convex set with diameter  $D$  and  $(\ell_t)_{t=1}^T$  a sequence of loss functions with each  $\ell_t : K \rightarrow \mathbb{R}_+$   $C$ -bounded and  $L$ -Lipschitz. Let  $R$  be a  $\gamma$ -self-concordant barrier for  $K$ . Then, for  $k \geq 12Cd$ , the regret of OPTIMISTIC BCO can be bounded as follows:  $Ck^2 Cd^2 \sqrt{T} \text{RegT}(\text{OPTIMISTIC BCO}) \leq \sqrt{LT} + LDT + 2 + \log(1/\gamma) \cdot 2^2 (k+1) \cdot \frac{p}{\#p} \frac{1}{2} p^{48d} k + LT \cdot 2^2 D 3L + 2DLk + \dots$

In particular, for  $\gamma = T^{11/16} d$  for the regret of the algorithm:

$$\begin{aligned} &\frac{3}{8} \\ &\gamma = T \\ &\frac{5}{16} \frac{3}{8} \\ &d \\ &k = T^{1/8} d^{1/4}, \text{ the following guarantee holds} \\ &\gamma \leq T^{11/16} d^{3/8} \cdot \text{RegT}(\text{OPTIMISTIC BCO}) = O \end{aligned}$$

The above result is the first improvement on the regret of Lipschitz losses in terms of  $T$  since the original algorithm of Flaxman et al. [2005] that is realizable from a concrete algorithm as well as polynomial in both dimension and time (both computationally and in terms of regret). Theorem 2 ( $C_{1,1}$  Bound). Let  $K \subset \mathbb{R}^d$  be a convex set with diameter  $D$  and  $(\ell_t)_{t=1}^T$  a sequence of loss functions with each  $\ell_t : K \rightarrow \mathbb{R}_+$   $C$ -bounded,  $L$ -Lipschitz and  $H$ -smooth. Let  $R$  be a  $\gamma$ -selfconcordant barrier for  $K$ . Then, for  $k \geq 12d$ , the regret of OPTIMISTIC BCO can be bounded as follows:  $\text{RegT}(\text{OPTIMISTIC BCO}) \leq \sqrt{LT} + H^2 D^2 T \cdot \frac{p}{\#p} \frac{3L}{2} p^{48d} 1 + (T L + DHT)^2 kD + 2DL + p + \log(1/\gamma) + Ck + \dots$  In particular, for  $\gamma = T^{8/13} d$  for the regret of the algorithm:

$$\frac{5}{6}$$





$$\begin{aligned} & 2 \\ & k \\ & g^?t = \\ & 1 \sum_{i=0}^{k+1} gbt(i) \end{aligned}$$

Assume that we play  $y_t$  at every round. Then we have the structural estimate:  $\sum_{t=1}^T \sum_{k=0}^K E[fbt(x_t) fbt(x^{??})] + LT \sup_{t=1} E[kxt(i) x_t k_2] + E[g^?t \sum_{t=1}^T \sum_{i=0}^K i^2 [1, T], i^2 [0, kt]]$

While we use averaged smoothed losses as in [Dekel et al., 2015], the analysis in this lemma is actually somewhat different. Because Dekel et al. [2015] always assume that the loss functions are in  $C_{1,1}$ , they elect to use the following decomposition:  $fbt(x_t)$

$f^?t(x_t) + f^?t(x_t) f^?t(x^{??}) + f^?t(x^{??}) fbt(x^{??})$ . This is because they can relate  $f^?t(x) = \sum_{i=0}^{k+1} f^?t(i) fbt(i(x))$  to  $\sum_{i=0}^{k+1} f^?t(i) fbt(i(x))$  using the fact that the gradients are Lipschitz. Since the gradients of  $C$  functions are not Lipschitz, we cannot use the same analysis. Instead, we use the decomposition  $fbt(x_t)$

$$\begin{aligned} fbt(x^{??}) &= fbt(x_t) \\ fbt(x^{??}) &= fbt(x_t) \\ fbt(i(x_t)) &+ fbt(i(x_t)) \end{aligned}$$

The next lemma affirms that we do indeed get the improved predictable component of the average gradient.

$$\begin{aligned} & f^?t(x^{??}) + f^?t(x^{??}) \leq (k+1)^2 \\ & fbt(x^{??}). \end{aligned}$$

factor from predicting the

Lemma 7 (C<sub>0,1</sub> Algorithmic bound on the averaged losses). Let  $(f_t)_{t=1}^T$  be a sequence of loss functions, and assume that  $f_t : K \rightarrow \mathbb{R}$  is  $C$ -bounded and  $L$ -Lipschitz, where  $K \subseteq \mathbb{R}^d$ . Let  $P_T x^? = \arg\min_{x \in K} \sum_{t=1}^T f_t(x)$ , and let  $x^{??} = \arg\min_{x \in K} \sum_{t=1}^T \text{dist}(y_t, x)$ . Assume that we play according to the algorithm with  $\eta \leq 1/(2Cd)$ . Then we maintain the following guarantee:  $\sum_{t=1}^T \sum_{k=0}^K$

$$\begin{aligned} & E[g^?t \sum_{t=1}^T \sum_{i=0}^K i^2 [1, T], i^2 [0, kt]] \\ & \leq 2Cd \sum_{t=1}^T \sum_{i=0}^K i^2 [1, T] + R(x^{??}) \cdot 2(k+1)^2 \end{aligned}$$

So far, we have demonstrated a bound on the regret of the form:  $\sum_{k=0}^K \text{Reg}_T^{(A)} \leq \sum_{t=1}^T \sum_{k=0}^K i^2 [1, T] + LT \sup_{t=1} E[kxt(i) x_t k_2] +$

$$\begin{aligned} & 2Cd \sum_{t=1}^T \sum_{i=0}^K i^2 [1, T] + R(x^{??}) \cdot 2(k+1)^2 \end{aligned}$$

Thus, it remains to find a tight bound on  $\sup_{t=1}^T \sum_{i=0}^K i^2 [1, T] E[kxt(i) x_t k_2]$ , which measures the stability of the actions across the history that we average over. This result is similar to that of Dekel et al. [2015], except that we additionally need to account for the optimistic gradient prediction used. Lemma 8 (C<sub>0,1</sub> Algorithmic bound on the stability of actions). Let  $(f_t)_{t=1}^T$  be a sequence of loss functions, and assume that  $f_t : K \rightarrow \mathbb{R}$  is  $C$ -bounded and  $L$ -Lipschitz, where  $K \subseteq \mathbb{R}^d$ . Assume that we play according to the algorithm with  $\eta \leq 1/(2Cd)$ . Then the following estimate holds:  $\sum_{k=0}^K \sum_{t=1}^T \sum_{i=0}^K i^2 [1, T] E[kxt(i) x_t k_2] \leq 2\eta kD + 2DL + p \cdot k$  Proof. [of Theorem 1] Putting all the pieces

together from Lemmas 5, 6, 7, 8, shows that  $p \leq p^{1/2} \leq Ck^{2/3} T^{1/3} \Delta C_d$ .  $k \text{RegT}(A) \leq LT + LDT + 2 + R(x^*) + LT^{2/3} D^{1/3} 3L + 2DLk + \frac{1}{2}(k+1)^2$ . Since  $x^*$  is at least  $\frac{1}{2}$  away from the boundary, it follows from [Abernethy and Rakhlin, 2009] that  $R(x^*) \leq \log(1/\frac{1}{2})$ . Plugging in the stated quantities for  $\Delta$ ,  $k$ , and yields the result.  $\square$

#### C1,1 analysis

The analysis of the C1,1 regret bound is similar to the C0,1 case. The only difference is that we leverage the higher regularity of the losses to provide a more refined estimate on the cost of approximating  $f_t$  with  $\hat{f}_t$ . Apart from that, we will reuse the bounds derived in Lemmas 6, 7, and 8. The proof of the following lemma, along with that of Theorem 2, is provided in Appendix 7.3. Lemma 9 (C1,1 Structural bound on true losses in terms of smoothed losses). Let  $(f_t)_{t=1}^T$  be a sequence of loss functions, and assume that  $f_t : K \rightarrow \mathbb{R}$  is  $C$ -bounded,  $L$ -Lipschitz, and  $H$ -smooth, where  $K \subseteq \mathbb{R}^d$ . Denote  $d_{f_t}(x) = \mathbb{E}[f_t(x + At) + f_t(y_t) - f_t(x_t) - \langle \nabla f_t(x_t), At \rangle]$ .

For arbitrary  $At$ ,  $y_t$ , and  $u_t$ . Let  $x_t^* = \arg\min_{x \in K} f_t(x)$ , and let  $x_t^{**} = \arg\min_{y \in K, \text{dist}(y, \partial K) \geq k} f_t(y)$ . Assume that we play  $y_t$  at every round. Then the following structural estimate holds:  $T \text{RegT}(A) = \mathbb{E}[f_t(y_t) - f_t(x_t^*)] \leq LT + 2H^2 D^2 T + \mathbb{E}[d_{f_t}(x_t) - d_{f_t}(x_t^{**})]$ .  $\square$

6

$t=1$

#### Conclusion

We designed a computationally efficient algorithm for bandit convex optimization admitting state-of-the-art guarantees for C0,1 and C1,1 loss functions. This was achieved using the general and powerful technique of predicting predictable information re-use. The ideas we describe here are directly applicable to other areas of optimization, in particular stochastic optimization. Acknowledgements This work was partly funded by NSF CCF-1535987 and IIS-1618662 and NSF GRFP DGE-1342536.  $\square$

## 2 References

J. Abernethy and A. Rakhlin. Beating the adaptive bandit with high probability. In COLT, 2009. J. Abernethy, E. Hazan, and A. Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. In COLT, pages 263–274, 2008. A. Agarwal, O. Dekel, and L. Xiao. Optimal algorithms for online convex optimization with multi-point bandit feedback. In COLT, pages 28–40, 2010. S. Bubeck and R. Eldan. Multi-scale exploration of convex functions and bandit convex optimization. CoRR, abs/1507.06580, 2015. S. Bubeck, O. Dekel, T. Koren, and Y. Peres. Bandit convex optimization:  $\sqrt{T}$  regret in one dimension. CoRR, abs/1502.06398, 2015. S. Bubeck, R. Eldan, and Y. T. Lee. Kernel-based methods for bandit convex optimization. CoRR, abs/1607.03084, 2016. V. Dani, T. P. Hayes, and S. M. Kakade. Stochastic linear optimization under bandit feedback. A. Defazio, F. Bach, and

S. Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In NIPS, pages 1646?1654, 2014. O. Dekel, R. Eldan, and T. Koren. Bandit smooth convex optimization: Improving the bias-variance tradeoff. In NIPS, pages 2908?2916, 2015. A. D. Flaxman, A. T. Kalai, and H. B. McMahan. Online convex optimization in the bandit setting: Gradient descent without a gradient. In SODA, pages 385?394, 2005. E. Hazan and Y. Li. An optimal algorithm for bandit convex optimization. CoRR, abs/1603.04350, 2016. R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In NIPS, pages 315?323, 2013. R. D. Kleinberg. Nearly tight bounds for the continuum-armed bandit problem. In Advances in Neural Information Processing Systems, pages 697?704, 2004. Y. Nesterov. Introductory Lectures on Convex Optimization: A Basic Course. Springer, New York, NY, USA, 2004. Y. Nesterov and A. Nemirovskii. Interior-point Polynomial Algorithms in Convex Programming. Studies in Applied Mathematics. Society for Industrial and Applied Mathematics, 1994. ISBN 9781611970791. A. Rakhlin and K. Sridharan. Online learning with predictable sequences. In COLT, pages 993?1019, 2013. A. Saha and A. Tewari. Improved regret guarantees for online smooth convex optimization with bandit feedback. In AISTATS, pages 636?642, 2011. M. W. Schmidt, N. L. Roux, and F. R. Bach. Minimizing finite sums with the stochastic average gradient. CoRR, abs/1309.2388, 2013.