

---

# Efficient Online Linear Optimization with Approximation Algorithms

---

Dan Garber

Technion - Israel Institute of Technology  
dangar@technion.ac.il

## Abstract

We revisit the problem of *online linear optimization* in case the set of feasible actions is accessible through an approximated linear optimization oracle with a factor  $\alpha$  multiplicative approximation guarantee. This setting is in particular interesting since it captures natural online extensions of well-studied *offline* linear optimization problems which are NP-hard, yet admit efficient approximation algorithms. The goal here is to minimize the  $\alpha$ -*regret* which is the natural extension of the standard *regret* in *online learning* to this setting. We present new algorithms with significantly improved oracle complexity for both the full information and bandit variants of the problem. Mainly, for both variants, we present  $\alpha$ -regret bounds of  $O(T^{-1/3})$ , where  $T$  is the number of prediction rounds, using only  $O(\log(T))$  calls to the approximation oracle per iteration, on average. These are the first results to obtain both average oracle complexity of  $O(\log(T))$  (or even poly-logarithmic in  $T$ ) and  $\alpha$ -regret bound  $O(T^{-c})$  for a constant  $c > 0$ , for both variants.

## 1 Introduction

In this paper we revisit the problem of *Online Linear Optimization* (OLO) [14], which is a specialized case of *Online Convex Optimization* (OCO) [12] with linear loss functions, in case the feasible set of actions is accessible through an oracle for approximated linear optimization with a multiplicative approximation error guarantee. In the standard setting of OLO, a decision maker is repeatedly required to choose an action, a vector in some fixed feasible set in  $\mathbb{R}^d$ . After choosing his action, the decision maker incurs loss (or payoff) given by the inner product between his selected vector and a vector chosen by an adversary. This game between the decision maker and the adversary then repeats itself. In the *full information* variant of the problem, after the decision maker receives his loss (payoff) on a certain round, he gets to observe the vector chosen by the adversary. In the *bandit* version of the problem, the decision maker only observes his loss (payoff) and does not get to observe the adversary's vector. The standard goal of the decision maker in OLO is to minimize a quantity known as *regret*, which measures the difference between the average loss of the decision maker on a game of  $T$  consecutive rounds (where  $T$  is fixed and known in advance), and the average loss of the best feasible action in hindsight (i.e., chosen with knowledge of all actions of the adversary throughout the  $T$  rounds) (in case of payoffs this difference is reversed). The main concern when designing algorithms for choosing the actions of the decision maker, is guaranteeing that the regret goes to zero as the length of the game  $T$  increases, as fast as possible (i.e., the rate of the regret in terms of  $T$ ). It should be noted that in this paper we focus on the case in which the adversary is *oblivious* (a.k.a. *non-adaptive*), which means the adversary chooses his entire sequence of actions for the  $T$  rounds beforehand.

While there exist well known algorithms for choosing the decision maker's actions which guarantee optimal regret bounds in  $T$ , such as the celebrated *Follow the Perturbed Leader* (FPL) and *Online Gradient Descent* (OGD) algorithms [14, 17, 12], efficient implementation of these algorithms hinges

on the ability to efficiently solve certain convex optimization problems (e.g., linear minimization for FPL or Euclidean projection for OGD) over the feasible set (or the convex hull of feasible points). However, when the feasible set corresponds for instance to the set of all possible solutions to some NP-Hard optimization problem, no such efficient implementations are known (or even widely believed to exist), and thus these celebrated regret-minimizing procedures cannot be efficiently applied. Luckily, many NP-Hard linear optimization problems (i.e., the objective function to either minimize or maximize is linear) admit efficient approximation algorithms with a multiplicative approximation guarantee. Some examples include MAX-CUT (factor 0.87856 approximation due to [9]), METRIC TSP (factor 1.5 approximation due to [6]), MINIMUM WEIGHTED VERTEX COVER (factor 2 approximation [4]), and WEIGHTED SET COVER (factor  $(\log n + 1)$  approximation due to [7]). It is thus natural to ask whether an efficient factor  $\alpha$  approximation algorithm for an NP-Hard *offline* linear optimization problem could be used to construct, in a generic way, an efficient algorithm for the *online* version of the problem. Note that in this case, even efficiently computing the best fixed action in hindsight is not possible, and thus, minimizing regret via an efficient algorithm does not seem likely (given an approximation algorithm we can however compute in hindsight a decision that corresponds to at most (at least)  $\alpha$  times the average loss (payoff) of the best fixed decision in hindsight).

In their paper [13], Kakade, Kalai and Ligett were the first to address this question in a fully generic way. They showed that using only an  $\alpha$ -approximation oracle for the set of feasible actions, it is possible, at a high level, to construct an online algorithm which achieves vanishing (expected)  $\alpha$ -regret, which is the difference between the average loss of the decision maker and  $\alpha$  times the average loss of the best fixed point in hindsight (for loss minimization problems and  $\alpha \geq 1$ ; a corresponding definition exists for payoff maximization problems and  $\alpha < 1$ ). Concretely, [13] showed that one can guarantee  $O(T^{-1/2})$  expected  $\alpha$ -regret in the full-information setting, which is optimal, and  $O(T^{-1/3})$  in the bandit setting under the additional assumption of the availability of a *Barycentric Spanner* (which we discuss in the sequel).

While the algorithm in [13] achieves an optimal  $\alpha$ -regret bound (in terms of  $T$ ) for the full information setting, in terms of computational complexity, the algorithm requires, in worst case, to perform on each round  $O(T)$  calls to the approximation oracle, which might be prohibitive and render the algorithm inefficient, since as discussed, in general,  $T$  is assumed to grow to infinity and thus the dependence of the runtime on  $T$  is of primary interest. Similarly, their algorithm for the bandit setting requires  $O(T^{2/3})$  calls to the approximation oracle per iteration.

The main contribution of our work is in providing new low  $\alpha$ -regret algorithms for the full information and bandit settings with significantly improved oracle complexities. A detailed comparison with [13] is given in Table 1. Concretely, for the full-information setting, we show it is possible to achieve  $O(T^{-1/3})$  expected  $\alpha$ -regret using only  $O(\log(T))$  calls to the approximation oracle per iteration, on average, which significantly improves over the  $O(T)$  bound of [13]<sup>1</sup>. We also show a bound of  $O(T^{-1/2})$  on the expected  $\alpha$ -regret (which is optimal) using only  $O(\sqrt{T} \log(T))$  calls to the oracle per iteration, on average, which gives nearly quadratic improvement over [13]. In the bandit setting we show it is possible to obtain a  $O(T^{-1/3})$  bound on the expected  $\alpha$ -regret (same as in [13]) using only  $O(\log(T))$  calls to the oracle per iteration, on average, under the same assumption on the availability of a *Barycentric Spanner* (BS). It is important to note that while there exist algorithms for OLO with bandit feedback which guarantee  $\tilde{O}(T^{-1/2})$  expected regret [1, 11] (where the  $\tilde{O}(\cdot)$  hides poly-logarithmic factors in  $T$ ), these require on each iteration to either solve to arbitrarily small accuracy a convex optimization problem over the feasible set [1], or sample a point from the feasible set according to a specified distribution [11], both of which cannot be implemented efficiently in our setting. On the other-hand, as we formally show in the sequel, at a high level, using a BS (originally introduced in [2]) simply requires to find a single set of  $d$  points from the feasible set which span the entire space  $\mathbb{R}^d$  (assuming this is possible, otherwise the set could be mapped to a lower dimensional space). The process of finding these vectors can be viewed as a preprocessing step and thus can be carried out offline. Moreover, as discussed in [13], for many NP-Hard problems it is possible to compute a BS in polynomial time and thus even this preprocessing step is efficient. Importantly, [13] shows that the approximation oracle by itself is not strong enough to guarantee non-trivial  $\alpha$ -regret in the bandit setting, and hence this assumption on the availability of a BS seems reasonable. Since the

<sup>1</sup>as we show in the appendix, even if we relax the algorithm of [13] to only guarantee  $O(T^{-1/3})$   $\alpha$ -regret, it will still require  $O(T^{2/3})$  calls to the oracle per iteration, on average.

	full information		bandit information	
Reference	$\alpha$ – regret	oracle complexity	$\alpha$ – regret	oracle complexity
KKL [13]	$T^{-1/2}$	$T$	$T^{-1/3}$	$T^{2/3}$
This paper (Thm. 4.1, 4.2)	$T^{-1/3}$	$\log(T)$	$T^{-1/3}$	$\log(T)$
This paper (Thm. 4.1)	$T^{-1/2}$	$\sqrt{T} \log(T)$	-	-

Table 1: comparison of expected  $\alpha$  – regret bounds and average number of calls to the approximation oracle per iteration. In all bounds we give only the dependence on the length of the game  $T$  and omit all other dependencies which we treat as constants. In the bandit setting we report the *expected* number of calls to the oracle per iteration.

best general regret bound known using a BS is  $O(T^{-1/3})$ , the  $\alpha$ -regret bound of our bandit algorithm is the best achievable to date via an efficient algorithm.

Technically, the main challenge in the considered setting is that as discussed, we cannot readily apply standard tools such as FPL and OGD. At a high level, in [13] it was shown that it is possible to apply the OGD method by replacing the exact projection step of OGD with an iterative algorithm which finds an *infeasible* point, but one that both satisfies the projection property required by OGD and is dominated by a convex combination of feasible points for every relevant linear loss (payoff) function. Unfortunately, in worst case, the number of queries to the approximation oracle required by this so-called projection algorithm per iteration is linear in  $T$ . While our online algorithms are also based on an application of OGD, our approach to computing the so-called projections is drastically different than [13], and is based on a coupling of two *cutting plane methods*, one that is based on the Ellipsoid method, and the other that resembles Gradient Descent. This approach might be of independent interest and might prove useful to similar problems.

### 1.1 Additional related work

Kalai and Vempala [14] showed that approximation algorithms which have *point-wise approximation guarantee*, such as the celebrated MAX-CUT algorithm of [9], could be used to instantiate their *Follow the Perturbed Leader* framework to achieve low  $\alpha$ -regret. However this construction is far from generic and requires the oracle to satisfy additional non-trivial conditions. This approach was also used in [3]. In [14] it was also shown that FPL could be instantiated with a FPTAS to achieve low  $\alpha$ -regret, however the approximation factor in the FPTAS needs to be set to roughly  $(1 + O(T^{-1/2}))$ , which may result in prohibitive running times even if a FPTAS for the underlying problem is available. Similarly, in [8] it was shown that if the approximation algorithm is based on solving a convex relaxation of the original, possibly NP-Hard, problem, this additional structure can be used with the FPL framework to achieve low  $\alpha$ -regret efficiently. To conclude all of the latter works consider specialized cases in which the approximation oracle satisfies additional non-trivial assumptions beyond its approximation guarantee, whereas here, similarly to [13], we will be interested in a generic as possible conversion from the offline problem to the online one, without imposing additional structure on the offline oracle.

## 2 Preliminaries

### 2.1 Online linear optimization with approximation oracles

Let  $\mathcal{K}, \mathcal{F}$  be compact sets of points in  $\mathbb{R}_+^d$  (non-negative orthant in  $\mathbb{R}^d$ ) such that  $\max_{\mathbf{x} \in \mathcal{K}} \|\mathbf{x}\| \leq R$ ,  $\max_{\mathbf{f} \in \mathcal{F}} \|\mathbf{f}\| \leq F$ , for some  $R > 0, F > 0$  (throughout this work we let  $\|\cdot\|$  denote the standard Euclidean norm), and for all  $\mathbf{x} \in \mathcal{K}, \mathbf{f} \in \mathcal{F}$  it holds that  $C \geq \mathbf{x} \cdot \mathbf{f} \geq 0$ , for some  $C > 0$ .

We assume  $\mathcal{K}$  is accessible through an approximated linear optimization oracle  $\mathcal{O}_{\mathcal{K}} : \mathbb{R}_+^d \rightarrow \mathcal{K}$  with parameter  $\alpha > 0$  such that:

$$\forall \mathbf{c} \in \mathbb{R}_+^d : \quad \mathcal{O}_{\mathcal{K}}(\mathbf{c}) \in \mathcal{K} \quad \text{and} \quad \begin{cases} \mathcal{O}_{\mathcal{K}}(\mathbf{c}) \cdot \mathbf{c} \leq \alpha \min_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot \mathbf{c} & \text{if } \alpha \geq 1; \\ \mathcal{O}_{\mathcal{K}}(\mathbf{c}) \cdot \mathbf{c} \geq \alpha \max_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot \mathbf{c} & \text{if } \alpha < 1. \end{cases}$$

Here  $\mathcal{K}$  is the feasible set of actions for the player, and  $\mathcal{F}$  is the set of all possible loss/payoff vectors<sup>2</sup>.

<sup>2</sup>we note that both of our assumptions that  $\mathcal{K} \subset \mathbb{R}_+^d, \mathcal{F} \subset \mathbb{R}_+^d$  and that the oracle takes inputs from  $\mathbb{R}_+^d$  are made for ease of presentation and clarity, and since these naturally hold for many NP-Hard optimization problem that are relevant to our setting. Nevertheless, these assumptions could be easily generalized as done in [13].

Since naturally a factor  $\alpha > 1$  for the approximation oracle is reasonable only for loss minimization problems, and a value  $\alpha < 1$  is reasonable for payoff maximization problems, throughout this work it will be convenient to use the value of  $\alpha$  to differentiate between minimization problems and maximization problems.

Given a sequence of linear loss/payoff functions  $\{\mathbf{f}_1, \dots, \mathbf{f}_T\} \in \mathcal{F}^T$  and a sequence of feasible points  $\{\mathbf{x}_1, \dots, \mathbf{x}_T\} \in \mathcal{K}^T$ , we define the  $\alpha$ -regret of the sequence  $\{\mathbf{x}_t\}_{t \in [T]}$  with respect to the sequence  $\{\mathbf{f}_t\}_{t \in [T]}$  as

$$\alpha - \text{regret}(\{\mathbf{x}_t, \mathbf{f}_t\}_{t \in [T]}) := \begin{cases} \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \cdot \mathbf{f}_t - \alpha \cdot \min_{\mathbf{x} \in \mathcal{K}} \frac{1}{T} \sum_{t=1}^T \mathbf{x} \cdot \mathbf{f}_t & \text{if } \alpha \geq 1; \\ \alpha \cdot \max_{\mathbf{x} \in \mathcal{K}} \frac{1}{T} \sum_{t=1}^T \mathbf{x} \cdot \mathbf{f}_t - \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \cdot \mathbf{f}_t & \text{if } \alpha < 1. \end{cases} \quad (1)$$

When the sequences  $\{\mathbf{x}_t\}_{t \in [T]}$ ,  $\{\mathbf{f}_t\}_{t \in [T]}$  are obvious from context we will simply write  $\alpha$ -regret without stating these sequences. Also, when the sequence  $\{\mathbf{x}_t\}_{t \in [T]}$  is randomized we will use  $\mathbb{E}[\alpha - \text{regret}]$  to denote the expected  $\alpha$ -regret.

### 2.1.1 Online linear optimization with full information

In OLO with full information, we consider a repeated game of  $T$  prediction rounds, for a fixed  $T$ , where on each round  $t$ , the decision maker is required to choose a feasible action  $\mathbf{x}_t \in \mathcal{K}$ . After committing to his choice, a linear loss function  $\mathbf{f}_t \in \mathcal{F}$  is revealed, and the decision maker incurs loss of  $\mathbf{x}_t \cdot \mathbf{f}_t$ . In the payoff version, the decision maker incurs payoff of  $\mathbf{x}_t \cdot \mathbf{f}_t$ . The game then continues to the next round. The overall goal of the decision maker is to guarantee that  $\alpha - \text{regret}(\{\mathbf{x}_t, \mathbf{f}_t\}_{t \in [T]}) = O(T^{-c})$  for some  $c > 0$ , at least in expectation (in fact using randomization is mandatory since  $\mathcal{K}$  need not be convex). Here we assume that the adversary is *oblivious* (aka *non-adaptive*), i.e., the sequence of losses/payoffs  $\mathbf{f}_1, \dots, \mathbf{f}_T$  is chosen in advance (before the first round), and does not depend on the actions of the decision maker.

### 2.1.2 Bandit feedback

The bandit version of the problem is identical to the full information setting with one crucial difference: on each round  $t$ , after making his choice, the decision maker does not observe the vector  $\mathbf{f}_t$ , but only the value of his loss/payoff, given by  $\mathbf{x}_t \cdot \mathbf{f}_t$ .

## 2.2 Additional notation

For any two sets  $\mathcal{S}, \mathcal{K} \subset \mathbb{R}^d$  and a scalar  $\beta \in \mathbb{R}$  we define the sets  $\mathcal{S} + \mathcal{K} := \{\mathbf{x} + \mathbf{y} \mid \mathbf{x} \in \mathcal{S}, \mathbf{y} \in \mathcal{K}\}$ ,  $\beta\mathcal{S} := \{\beta\mathbf{x} \mid \mathbf{x} \in \mathcal{S}\}$ . We also denote by  $\text{CH}(\mathcal{K})$  the convex-hull of all points in a set  $\mathcal{K}$ . For a convex and compact set  $\mathcal{S} \subset \mathbb{R}^d$  and a point  $\mathbf{x} \in \mathbb{R}^d$  we define  $\text{dist}(\mathbf{x}, \mathcal{S}) := \min_{\mathbf{z} \in \mathcal{S}} \|\mathbf{z} - \mathbf{x}\|$ . We let  $\mathcal{B}(\mathbf{c}, r)$  denote the Euclidean ball or radius  $r$  centered in  $\mathbf{c}$ .

## 2.3 Basic algorithmic tools

We now briefly describe two very basic ideas that are essential for constructing our algorithms, namely the *extended approximation oracle* and the *online gradient descent without feasibility* method. These were already suggested in [13] to obtain their low  $\alpha$ -regret algorithms. We note that in the appendix we describe in more detail the approach of [13] and discuss its shortcomings in obtaining oracle-efficient algorithms.

### 2.3.1 The extended approximation oracle

As discussed, a key difficulty of our setting that prevents us from directly applying well studied algorithms for OLO, is that essentially all standard algorithms require to exactly solve (or up to arbitrarily small error) some linear/convex optimization problem over the convexification of the feasible set  $\text{CH}(\mathcal{K})$ . However, not only that our approximation oracle  $\mathcal{O}_{\mathcal{K}}(\cdot)$  cannot perform exact minimization, even for  $\alpha = 1$  it is applicable only with inputs in  $\mathbb{R}_+^d$ , and hence cannot optimize in all directions. A natural approach, suggested in [13], to overcome the approximation error of the oracle  $\mathcal{O}_{\mathcal{K}}(\cdot)$ , is to consider optimization with respect to the convex set  $\text{CH}(\alpha\mathcal{K})$  (i.e. convex hull of all points in  $\mathcal{K}$  scaled by a factor of  $\alpha$ ) instead of  $\text{CH}(\mathcal{K})$ . Indeed, if we consider for instance the case  $\alpha \geq 1$ , it is straightforward to see that for any  $\mathbf{c} \in \mathbb{R}_+^d$ ,  $\mathcal{O}_{\mathcal{K}}(\mathbf{c}) \cdot \mathbf{c} \leq \alpha \min_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot \mathbf{c} =$

$\alpha \min_{\mathbf{x} \in \text{CH}(\mathcal{K})} \mathbf{x} \cdot \mathbf{c} = \min_{\mathbf{x} \in \text{CH}(\alpha\mathcal{K})} \mathbf{x} \cdot \mathbf{c}$ . Thus, in a certain sense,  $\mathcal{O}_{\mathcal{K}}(\cdot)$  can optimize with respect to  $\text{CH}(\alpha\mathcal{K})$  for all directions in  $\mathbb{R}_+^d$ , although the oracle returns points in the original set  $\mathcal{K}$ .

The following lemma shows that one can easily extend the oracle  $\mathcal{O}_{\mathcal{K}}(\cdot)$  to optimize with respect to all directions in  $\mathbb{R}^d$ .

**Lemma 2.1** (Extended approximation oracle). *Given  $\mathbf{c} \in \mathbb{R}^d$  write  $\mathbf{c} = \mathbf{c}^+ + \mathbf{c}^-$  where  $\mathbf{c}^+$  equals to  $\mathbf{c}$  on all non-negative coordinates of  $\mathbf{c}$  and zero everywhere else, and  $\mathbf{c}^-$  equals  $\mathbf{c}$  on all negative coordinates and zero everywhere else. The extended approximation oracle is a mapping  $\hat{\mathcal{O}}_{\mathcal{K}} : \mathbb{R}^d \rightarrow (\mathcal{K} + \mathcal{B}(0, (1 + \alpha)R), \mathcal{K})$  defined as:*

$$\hat{\mathcal{O}}_{\mathcal{K}}(\mathbf{c}) = (\mathbf{v}, \mathbf{s}) := \begin{cases} (\mathcal{O}_{\mathcal{K}}(\mathbf{c}^+) - \alpha R \bar{\mathbf{c}}^-, \mathcal{O}_{\mathcal{K}}(\mathbf{c}^+)) & \text{if } \alpha \geq 1; \\ (\mathcal{O}_{\mathcal{K}}(-\mathbf{c}^-) - R \bar{\mathbf{c}}^+, \mathcal{O}_{\mathcal{K}}(-\mathbf{c}^-)) & \text{if } \alpha < 1, \end{cases} \quad (2)$$

where for any vector  $\mathbf{v} \in \mathbb{R}^d$  we denote  $\bar{\mathbf{v}} = \mathbf{v}/\|\mathbf{v}\|$  if  $\|\mathbf{v}\| > 0$  and  $\bar{\mathbf{v}} = \mathbf{0}$  otherwise, and it satisfies the following three properties:

1.  $\mathbf{v} \cdot \mathbf{c} \leq \min_{\mathbf{x} \in \alpha\mathcal{K}} \mathbf{x} \cdot \mathbf{c}$
2.  $\forall \mathbf{f} \in \mathcal{F} : \mathbf{s} \cdot \mathbf{f} \leq \mathbf{v} \cdot \mathbf{f}$  if  $\alpha \geq 1$  and  $\mathbf{s} \cdot \mathbf{f} \geq \mathbf{v} \cdot \mathbf{f}$  if  $\alpha < 1$
3.  $\|\mathbf{v}\| \leq (\alpha + 2)R$

The proof is given in the appendix for completeness.

It is important to note that while the extended oracle provides solutions with values at least as low as any point in  $\text{CH}(\alpha\mathcal{K})$ , still in general the output point  $\mathbf{v}$  need not be in either  $\mathcal{K}$  or  $\text{CH}(\alpha\mathcal{K})$ , which means that it is not a feasible point to play in our OLO setting, nor does it allow us to optimize over  $\text{CH}(\alpha\mathcal{K})$ . This is why we also need the oracle to output the feasible point  $\mathbf{s} \in \mathcal{K}$  which dominates  $\mathbf{v}$  for any possible loss/payoff vector in  $\mathcal{F}$ . While we will use the outputs  $\mathbf{v}$  to solve a certain optimization problem involving  $\text{CH}(\alpha\mathcal{K})$ , this dominance relation will be used to convert the solutions to these optimization problems into feasible plays for our OLO algorithms.

### 2.3.2 Online gradient descent with and without feasibility

As in [13], our online algorithms will be based on the well known *Online Gradient Descent* method (OGD) for online convex optimization, originally due to [17]. For a sequence of loss vectors  $\{\mathbf{f}_1, \dots, \mathbf{f}_T\} \subset \mathbb{R}^d$  OGD produces a sequence of plays  $\{\mathbf{x}_1, \dots, \mathbf{x}_T\} \subset \mathcal{S}$ , for a convex and compact set  $\mathcal{S} \subset \mathbb{R}^d$  via the following updates:  $\forall t \geq 1 : \mathbf{y}_{t+1} \leftarrow \mathbf{x}_t - \eta \mathbf{f}_t$ ,  $\mathbf{x}_{t+1} \leftarrow \arg \min_{\mathbf{x} \in \mathcal{S}} \|\mathbf{x} - \mathbf{y}_{t+1}\|^2$ , where  $\mathbf{x}_1$  is initialized to some arbitrary point in  $\mathcal{S}$  and  $\eta$  is some pre-determined step-size. The obvious difficulty in applying OGD to online linear optimization over  $\mathcal{S} = \text{CH}(\alpha\mathcal{K})$  is the step of computing  $\mathbf{x}_{t+1}$  by projecting  $\mathbf{y}_{t+1}$  onto the feasible set  $\mathcal{S}$ , since as discussed, even with the extended approximation oracle, one cannot exactly optimize over  $\text{CH}(\alpha\mathcal{K})$ . Instead we will consider a variant of OGD which may produce infeasible points, i.e., outside of  $\mathcal{S}$ , but which guarantees low regret with respect to any point in  $\mathcal{S}$ . This algorithm, which we refer to as *online gradient descent without feasibility*, is given below (Algorithm 1).

---

#### Algorithm 1 Online Gradient Descent Without Feasibility

---

- 1: input: learning rate  $\eta > 0$
- 2:  $\mathbf{x}_1 \leftarrow$  some point in  $\mathcal{S}$
- 3: **for**  $t = 1 \dots T$  **do**
- 4:   play  $\mathbf{x}_t$  and receive loss/payoff vector  $\mathbf{f}_t \in \mathbb{R}^d$
- 5:    $\mathbf{y}_{t+1} \leftarrow \begin{cases} \mathbf{x}_t - \eta \mathbf{f}_t & \text{for losses} \\ \mathbf{x}_t + \eta \mathbf{f}_t & \text{for payoffs} \end{cases}$
- 6:   find  $\mathbf{x}_{t+1} \in \mathbb{R}^d$  such that

$$\forall \mathbf{z} \in \mathcal{S} : \quad \|\mathbf{z} - \mathbf{x}_{t+1}\|^2 \leq \|\mathbf{z} - \mathbf{y}_{t+1}\|^2 \quad (3)$$

- 7: **end for**
- 

**Lemma 2.2.** *[Online gradient descent without feasibility] Fix  $\eta > 0$ . Suppose Algorithm 1 is applied for  $T$  rounds and let  $\{\mathbf{f}_t\}_{t=1}^T \subset \mathbb{R}^d$  be the sequence of observed loss/payoff vectors, and let  $\{\mathbf{x}_t\}_{t=1}^T$*

be the sequence of points played by the algorithm. Then for any  $\mathbf{x} \in \mathcal{S}$  it holds that

$$\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \cdot \mathbf{f}_t - \frac{1}{T} \sum_{t=1}^T \mathbf{x} \cdot \mathbf{f}_t \leq \frac{1}{2T\eta} \|\mathbf{x}_1 - \mathbf{x}\|^2 + \frac{\eta}{2T} \sum_{t=1}^T \|\mathbf{f}_t\|^2 \quad \text{for losses;}$$

$$\frac{1}{T} \sum_{t=1}^T \mathbf{x} \cdot \mathbf{f}_t - \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \cdot \mathbf{f}_t \leq \frac{1}{2T\eta} \|\mathbf{x}_1 - \mathbf{x}\|^2 + \frac{\eta}{2T} \sum_{t=1}^T \|\mathbf{f}_t\|^2 \quad \text{for payoffs.}$$

The proof is given in the appendix for completeness.

### 3 Oracle-efficient Computation of (infeasible) Projections onto $\text{CH}(\alpha\mathcal{K})$

In this section we detail our main technical tool for obtaining oracle-efficient online algorithms, i.e., our algorithm for computing projections, in the sense of Eq. (3), onto the convex set  $\text{CH}(\alpha\mathcal{K})$ . Before presenting our projection algorithm, Algorithm 2 and detailing its theoretical guarantees, we first present the main algorithmic building block in the algorithm, which is described in the following lemma. Lemma 3.1 shows that for any point  $\mathbf{x} \in \mathbb{R}^d$ , we can either find a near-by point  $\mathbf{p}$  which is a convex combination of points outputted by the extended approximation oracle (and hence,  $\mathbf{p}$  is dominated by a convex combination of feasible points in  $\mathcal{K}$  for any vector in  $\mathcal{F}$ , as discussed in Section 2.3.1), or we can find a separating hyperplane that separates  $\mathbf{x}$  from  $\text{CH}(\alpha\mathcal{K})$  with sufficiently large margin. We achieve this by running the well known Ellipsoid method [10, 5] in a very specialized way. This application of the Ellipsoid method is similar in spirit to those in [15, 16], which applied this idea to computing *correlated equilibrium* in games and *algorithmic mechanism design*, though the implementation details and the way in which we apply this technique are quite different.

The proof of the following lemma is given in the appendix.

**Lemma 3.1** (Separation-or-Decomposition via the Ellipsoid method). *Fix  $\mathbf{x} \in \mathbb{R}^d$ ,  $\epsilon \in (0, (\alpha + 2)R]$ , and a positive integer  $N \geq cd^2 \ln \left( \frac{(\alpha+1)R + \|\mathbf{x}\|}{\epsilon} \right)$ , where  $c$  is a positive universal constant. Consider an attempt to apply the Ellipsoid method for  $N$  iterations to the following feasibility problem:*

$$\text{find } \mathbf{w} \in \mathbb{R}^d \text{ such that: } \quad \forall \mathbf{z} \in \alpha\mathcal{K} : \quad (\mathbf{x} - \mathbf{z}) \cdot \mathbf{w} \geq \epsilon \quad \text{and} \quad \|\mathbf{w}\| \leq 1, \quad (4)$$

such that each iteration of the Ellipsoid method applies the following consecutive steps:

1.  $(\mathbf{v}, \mathbf{s}) \leftarrow \hat{\mathcal{O}}_{\mathcal{K}}(-\mathbf{w})$ , where  $\mathbf{w}$  is the current iterate. If  $(\mathbf{x} - \mathbf{v}) \cdot \mathbf{w} < \epsilon$ , use  $\mathbf{v} - \mathbf{x}$  as a separating hyperplane for the Ellipsoid method and continue to the next iteration
2. if  $\|\mathbf{w}\| > 1$ , use  $\mathbf{w}$  as a separating hyperplane for the Ellipsoid method and continue to the next iteration
3. otherwise ( $\|\mathbf{w}\| \leq 1$  and  $(\mathbf{x} - \mathbf{v}) \cdot \mathbf{w} \geq \epsilon$ ), declare Problem (4) feasible and return the vector  $\mathbf{w}$ .

Then, if the Ellipsoid method terminates declaring Problem 4 feasible, the returned vector  $\mathbf{w}$  is a feasible solution to Problem (4). Otherwise (the Ellipsoid method completes  $N$  iterations without declaring Problem (4) feasible), let  $(\mathbf{v}_1, \mathbf{s}_1), \dots, (\mathbf{v}_N, \mathbf{s}_N)$  be the outputs of the extended approximation oracle gathered throughout the run of the algorithm, and let  $(a_1, \dots, a_N)$  be an optimal solution to the following convex optimization problem:

$$\min_{(a_1, \dots, a_N)} \frac{1}{2} \left\| \sum_{i=1}^N a_i \mathbf{v}_i - \mathbf{x} \right\|^2 \quad \text{such that} \quad \forall i \in \{1, \dots, N\} : a_i \geq 0, \quad \sum_{i=1}^N a_i = 1. \quad (5)$$

Then the point  $\mathbf{p} = \sum_{i=1}^N a_i \mathbf{v}_i$  satisfies  $\|\mathbf{x} - \mathbf{p}\| \leq 3\epsilon$ .

We are now ready to present our algorithm for computing projections onto  $\text{CH}(\alpha\mathcal{K})$  (in the sense of Eq. (3)). Consider now an attempt to project a point  $\mathbf{y} \in \mathbb{R}^d$ , and note that in particular,  $\mathbf{y}$  itself is a valid projection (again, in the sense of Eq. (3)), however, in general, it is not a feasible point nor is it dominated by a convex combination of feasible points. When attempting to project  $\mathbf{y} \in \mathbb{R}^d$ , our algorithm continuously applies the *separation-or-decomposition* procedure described in Lemma 3.1.

In case the procedure returns a decomposition, then by Lemma 3.1, we have a point that is sufficiently close to  $\mathbf{y}$  and is dominated for any vector in  $\mathcal{F}$  by a convex combination (given explicitly) of feasible points in  $\mathcal{K}$ . Otherwise, the procedure returns a separating hyperplane which can be used to to “pull  $\mathbf{y}$  closer” to  $\text{CH}(\alpha\mathcal{K})$  in a way that the resulting point still satisfies the projection inequality given in Eq. (3), and the process then repeats itself. Since each time we obtain a hyperplane separating our current iterate from  $\text{CH}(\alpha\mathcal{K})$ , we pull the current iterate sufficiently towards  $\text{CH}(\alpha\mathcal{K})$ , this process must terminate. Lemma 3.2 gives exact bounds on the performance of the algorithm.

---

**Algorithm 2** (infeasible) Projection onto  $\text{CH}(\alpha\mathcal{K})$

---

```

1: input: point  $\mathbf{y} \in \mathbb{R}^d$ , tolerance  $\epsilon > 0$ 
2:  $\tilde{\mathbf{y}} \leftarrow \mathbf{y} / \max\{1, \|\mathbf{y}\|/(\alpha R)\}$ 
3: for  $t = 1 \dots$  do
4:   call the SEPARATION-OR-DECOMPOSTION procedure (Lemma 3.1) with parameters  $(\tilde{\mathbf{y}}, \epsilon)$ 
5:   if the procedure outputs a separating hyperplane  $\mathbf{w}$  then
6:      $\tilde{\mathbf{y}} \leftarrow \tilde{\mathbf{y}} - \epsilon \mathbf{w}$ 
7:   else
8:     let  $(a_1, \dots, a_N), \{(\mathbf{v}_1, \mathbf{s}_1), \dots, (\mathbf{v}_N, \mathbf{s}_N)\}$  be the decomposition returned
9:     return  $\tilde{\mathbf{y}}, (a_1, \dots, a_N), \{(\mathbf{v}_1, \mathbf{s}_1), \dots, (\mathbf{v}_N, \mathbf{s}_N)\}$ 
10:  end if
11: end for

```

---

**Lemma 3.2.** Fix  $\mathbf{y} \in \mathbb{R}^d$  and  $\epsilon \in (0, (\alpha + 2)R]$ . Algorithm 2 terminates after at most  $\lceil \alpha^2 R^2 / \epsilon^2 \rceil$  iterations, returning a point  $\tilde{\mathbf{y}} \in \mathbb{R}^d$ , a distribution  $(a_1, \dots, a_N)$  and a set  $\{(\mathbf{v}_1, \mathbf{s}_1), \dots, (\mathbf{v}_N, \mathbf{s}_N)\}$  outputted by the extended approximation oracle, where  $N$  is as defined in Lemma 3.1, such that

$$1. \quad \forall \mathbf{z} \in \text{CH}(\alpha\mathcal{K}) : \quad \|\tilde{\mathbf{y}} - \mathbf{z}\|^2 \leq \|\mathbf{y} - \mathbf{z}\|^2, \quad 2. \quad \|\mathbf{p} - \tilde{\mathbf{y}}\| \leq 3\epsilon \quad \text{for } \mathbf{p} := \sum_{i \in [N]} a_i \mathbf{v}_i.$$

Moreover, if the **for** loop was entered a total number of  $k$  times, then the final value of  $\tilde{\mathbf{y}}$  satisfies

$$\text{dist}^2(\tilde{\mathbf{y}}, \text{CH}(\alpha\mathcal{K})) \leq \min\{2\alpha^2 R^2, \text{dist}^2(\mathbf{y}, \text{CH}(\alpha\mathcal{K})) - (k - 1)\epsilon^2\},$$

and the overall number of queries to the approximation oracle is  $O(kd^2 \ln((\alpha + 1)R/\epsilon))$ .

It is important to note that the worst case iteration bound in Lemma 3.2 does not seem so appealing for our purposes, since it depends polynomially on  $1/\epsilon$ , and in our online algorithms naturally we will need to take  $\epsilon = O(T^{-c})$  for some  $c > 0$ , which seems to contradict our goal of achieving poly-logarithmic in  $T$  oracle complexity, at least on average. However, as Lemma 3.2 shows, the more iterations Algorithm 2 performs, the closer it brings its final iterate to the set  $\text{CH}(\alpha\mathcal{K})$ . Thus, as we will show when analyzing the oracle complexity of our online algorithms, while a single call to Algorithm 2 can be expensive, when calling it sequentially, where each input is a small perturbation of the output of the previous call, the average number of iterations performed per such call cannot be too high.

## 4 Efficient Algorithms for the Full Information and Bandit Settings

We now turn to present our online algorithms for the full-information and bandit settings together with their regret bounds and oracle-complexity guarantees.

### 4.1 Algorithm for the full information setting

Our algorithm for the full-information setting, Algorithm 3, is given below.

**Theorem 4.1.** [Main Theorem] Fix  $\eta > 0, \epsilon \in (0, (\alpha + 2)R]$ . Suppose Algorithm 3 is applied for  $T$  rounds and let  $\{\mathbf{f}_t\}_{t=1}^T \subseteq \mathcal{F}$  be the sequence of observed loss/payoff vectors, and let  $\{\mathbf{s}_t\}_{t=1}^T$  be the sequence of points played by the algorithm. Then it holds that

$$\mathbb{E} [\alpha - \text{regret}(\{\mathbf{s}_t, \mathbf{f}_t\}_{t \in [T]})] \leq \alpha^2 R^2 T^{-1} \eta^{-1} + \eta F^2 / 2 + 3F\epsilon,$$

and the average number of calls to the approximation oracle of  $\mathcal{K}$  per iteration is upper bounded by

$$K(\eta, \epsilon) := O\left(\left(1 + (\eta\alpha RF + \eta^2 F^2) \epsilon^{-2}\right) d^2 \ln((\alpha + 1)R/\epsilon)\right).$$

---

**Algorithm 3** Online Gradient Descent with Infeasible Projections onto  $\text{CH}(\alpha\mathcal{K})$ 

---

```
1: input: learning rate  $\eta > 0$ , projection error parameter  $\epsilon > 0$ 
2:  $\mathbf{s}_1 \leftarrow$  some point in  $\mathcal{K}$ ,  $\tilde{\mathbf{y}}_1 \leftarrow \alpha\mathbf{s}_1$ 
3: for  $t = 1 \dots T$  do
4:   play  $\mathbf{s}_t$  and receive loss/payoff vector  $\mathbf{f}_t \in \mathcal{F}$ 
5:    $\mathbf{y}_{t+1} \leftarrow \begin{cases} \tilde{\mathbf{y}}_t - \eta\mathbf{f}_t & \text{if } \alpha \geq 1 \\ \tilde{\mathbf{y}}_t + \eta\mathbf{f}_t & \text{if } \alpha < 1 \end{cases}$ 
6:   call Algorithm 2 with inputs  $(\mathbf{y}_{t+1}, \epsilon)$  to obtain an approximated projection  $\tilde{\mathbf{y}}_{t+1}$ , a distribution  $(a_1, \dots, a_N)$  and  $\{(\mathbf{v}_1, \mathbf{s}_1), \dots, (\mathbf{v}_N, \mathbf{s}_N)\} \subseteq \mathbb{R}^d \times \mathcal{K}$ , for some  $N \in \mathbb{N}$ .
7:   sample  $\mathbf{s}_{t+1} \in \{\mathbf{s}_1, \dots, \mathbf{s}_N\}$  according to distribution  $(a_1, \dots, a_N)$ 
8: end for
```

---

In particular, setting  $\eta = \alpha RT^{-2/3}/F$ ,  $\epsilon = \alpha RT^{-1/3}$  gives  $\mathbb{E}[\alpha - \text{regret}] = O(\alpha RFT^{-1/3})$ ,  $K = O(d^2 \ln(\frac{\alpha+1}{\alpha}T))$ . Alternatively, setting  $\eta = \alpha RT^{-1/2}/F$ ,  $\epsilon = \alpha RT^{-1/2}$  gives  $\mathbb{E}[\alpha - \text{regret}] = O(\alpha RFT^{-1/2})$ ,  $K = O(\sqrt{T}d^2 \ln(\frac{\alpha+1}{\alpha}T))$ .

The proof is given in the appendix.

## 4.2 Algorithm for the bandit information setting

Our algorithm for the bandit setting follows from a very well known reduction from the bandit setting to the full information setting, also applied in the bandit algorithm of [13]. The algorithm simply simulates the full information algorithm, Algorithm 3, by providing it with estimated loss/payoff vectors  $\hat{\mathbf{f}}_1, \dots, \hat{\mathbf{f}}_T$  instead of the true vectors  $\mathbf{f}_1, \dots, \mathbf{f}_T$  which are not available in the bandit setting. This reduction is based on the use of a *Barycentric Spanner* (defined next) for the feasible set  $\mathcal{K}$ . As standard, we assume the points in  $\mathcal{K}$  span the entire space  $\mathbb{R}^d$ , otherwise we can reformulate the problem in a lower-dimensional space, in which this assumption holds.

**Definition 4.1** (Barycentric Spanner<sup>3</sup>). We say that a set of  $d$  vectors  $\{\mathbf{q}_1, \dots, \mathbf{q}_d\} \subset \mathbb{R}^d$  is a *Barycentric Spanner* with parameter  $\beta > 0$  for a set  $\mathcal{S} \subset \mathbb{R}^d$ , denoted by  $\beta\text{-BS}(\mathcal{S})$ , if it holds that  $\{\mathbf{q}_1, \dots, \mathbf{q}_d\} \subset \mathcal{S}$ , and the matrix  $\mathbf{Q} := \sum_{i=1}^d \mathbf{q}_i \mathbf{q}_i^\top$  is not singular and  $\max_{i \in [d]} \|\mathbf{Q}^{-1} \mathbf{q}_i\| \leq \beta$ .

Importantly, as discussed in [13], the assumption on the availability of such a set  $\beta\text{-BS}(\mathcal{K})$  seems reasonable, since i) for many sets that correspond to the set of all possible solutions to some well-studied NP-Hard optimization problem, one can still construct in  $\text{poly}(d)$  time a barycentric spanner with  $\beta = \text{poly}(d)$ , ii)  $\beta\text{-BS}(\mathcal{K})$  needs to be constructed only once and then stored in memory (overall  $d$  vectors in  $\mathbb{R}^d$ ), and hence its construction can be viewed as a pre-processing step, and iii) as illustrated in [13], without further assumptions, the approximation oracle by itself is not sufficient to guarantee nontrivial regret bounds in the bandit setting.

The algorithm and the proof of the following theorem are given in the appendix.

**Theorem 4.2.** Fix  $\eta > 0, \epsilon \in (0, (\alpha + 2)R], \gamma \in (0, 1)$ . Suppose Algorithm 5 is applied for  $T$  rounds and let  $\{\mathbf{f}_t\}_{t=1}^T \subseteq \mathcal{F}$  be the sequence of observed loss/payoff vectors, and let  $\{\hat{\mathbf{s}}_t\}_{t=1}^T$  be the sequence of points played by the algorithm. Then it holds that

$$\mathbb{E}[\alpha - \text{regret}(\{(\hat{\mathbf{s}}_t, \mathbf{f}_t)\}_{t \in [T]})] \leq \alpha^2 R^2 \eta^{-1} T^{-1} + \eta d^2 C^2 \beta^2 \gamma^{-1} / 2 + 3\epsilon F + \gamma C,$$

and the expected number of calls to the approximation oracle of  $\mathcal{K}$  per iteration is upper bounded by

$$\mathbb{E}[K(\eta, \epsilon, \gamma)] := O\left((1 + (\eta\alpha\beta d C R + (\eta d C \beta)^2 / \gamma) \epsilon^{-2}) d^2 \ln((\alpha + 1)R/\epsilon)\right).$$

In particular, setting  $\eta = \frac{\alpha R}{\beta d C} T^{-2/3}$ ,  $\epsilon = \alpha R T^{-1/3}$ ,  $\gamma = T^{-1/3}$  gives  $\mathbb{E}[\alpha - \text{regret}] = O((\alpha\beta d C R + \alpha R F + C) T^{-1/3})$ ,  $\mathbb{E}[K] = O(d^2 \ln(\frac{\alpha+1}{\alpha}T))$ .

---

<sup>3</sup>this definition is somewhat different than the classical one given in [2], however it is equivalent to a  $C$ -approximate barycentric spanner [2], with an appropriately chosen constant  $C(\beta)$ .



## 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] Baruch Awerbuch and Robert D Kleinberg. Adaptive routing with end-to-end feedback: Distributed learning and geometric approaches. In *Proceedings of the thirty-sixth annual ACM symposium on Theory of computing*, pages 45–53. ACM, 2004.
- [3] Maria-Florina Balcan and Avrim Blum. Approximation algorithms and online mechanisms for item pricing. In *Proceedings of the 7th ACM Conference on Electronic Commerce*, pages 29–35. ACM, 2006.
- [4] Reuven Bar-Yehuda and Shimon Even. A linear-time approximation algorithm for the weighted vertex cover problem. *Journal of Algorithms*, 2(2):198–203, 1981.
- [5] Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- [6] Nicos Christofides. Worst-case analysis of a new heuristic for the travelling salesman problem. Technical report, DTIC Document, 1976.
- [7] V. Chvatal. A greedy heuristic for the set-covering problem. *Mathematics of Operations Research*, 4(3):233–235, 1979.
- [8] Takahiro Fujita, Kohei Hatano, and Eiji Takimoto. Combinatorial online prediction via metarounding. In *ALT*, pages 68–82. Springer, 2013.
- [9] Michel X Goemans and David P Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *Journal of the ACM (JACM)*, 42(6):1115–1145, 1995.
- [10] M. Grötschel, L. Lovász, and A. Schrijver. The ellipsoid method and its consequences in combinatorial optimization. *Combinatorica*, 1(2):169–197, 1981.
- [11] Elad Hazan, Zohar Shay Karnin, and Raghu Meka. Volumetric spanners: an efficient exploration basis for learning. In *COLT*, volume 35, pages 408–422, 2014.
- [12] Elad Hazan and Haipeng Luo. Variance-reduced and projection-free stochastic optimization. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, pages 1263–1271, 2016.
- [13] Sham M. Kakade, Adam Tauman Kalai, and Katrina Ligett. Playing games with approximation algorithms. *SIAM J. Comput.*, 39(3):1088–1106, 2009.
- [14] Adam Kalai and Santosh Vempala. Efficient algorithms for online decision problems. *Journal of Computer and System Sciences*, 71(3):291–307, 2005.
- [15] Christos H Papadimitriou and Tim Roughgarden. Computing correlated equilibria in multi-player games. *Journal of the ACM (JACM)*, 55(3):14, 2008.
- [16] S Matthew Weinberg. *Algorithms for strategic agents*. PhD thesis, Massachusetts Institute of Technology, 2014.
- [17] Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Machine Learning, Proceedings of the Twentieth International Conference (ICML 2003), August 21-24, 2003, Washington, DC, USA*, pages 928–936, 2003.

## A The KKL approach

We now briefly describe how [13] use the extended approximation oracle and the online gradient descent without feasibility approach to construct their low  $\alpha$ -regret algorithm for the full information setting, and point out the limitation of this approach to obtaining low oracle complexity.

Consider some iteration  $t$  of Algorithm 1 and let  $\mathbf{y}_{t+1}$  be the newly computed point. Let  $(\mathbf{x}, \mathbf{s}) \in \mathbb{R}^d \times \mathcal{K}$  be such that  $\forall \mathbf{f} \in \mathcal{F} : \mathbf{x} \cdot \mathbf{f} \geq \mathbf{s} \cdot \mathbf{f}$  (e.g., take  $\mathbf{x} = \mathbf{s}$ ), and let  $(\mathbf{v}', \mathbf{s}') \leftarrow \hat{\mathcal{O}}_{\mathcal{K}}(\mathbf{x} - \mathbf{y}_{t+1})$ . We have the following simple lemma.

**Lemma A.1.** *Fix  $\epsilon \in (0, 3(\alpha+2)^2 R^2]$  and suppose that  $\mathbf{x} \in \mathcal{B}(0, (\alpha+2)R)$ . If  $(\mathbf{x} - \mathbf{y}_{t+1}) \cdot (\mathbf{x} - \mathbf{v}') \leq \epsilon$ , then setting  $\mathbf{x}_{t+1} \leftarrow \mathbf{x}$  gives*

$$\forall \mathbf{z} \in CH(\alpha\mathcal{K}) : \quad \|\mathbf{z} - \mathbf{x}_{t+1}\|^2 \leq \|\mathbf{z} - \mathbf{y}_{t+1}\|^2 + 2\epsilon.$$

*Otherwise, setting  $\mathbf{x}' \leftarrow (1 - \lambda)\mathbf{x} + \lambda\mathbf{v}'$ , for appropriately chosen  $\lambda \in (0, 1)$ , guarantees that*

$$\|\mathbf{x}' - \mathbf{y}_{t+1}\|^2 \leq \|\mathbf{x} - \mathbf{y}_{t+1}\|^2 - \Omega(\epsilon^2),$$

*and*

$$\forall \mathbf{f} \in \mathcal{F} : \quad ((1 - \lambda)\mathbf{s} + \lambda\mathbf{s}') \cdot \mathbf{f} \leq \mathbf{x}' \cdot \mathbf{f}.$$

*Proof.* To prove the first part of the lemma, suppose that  $\mathbf{x}$  satisfies that  $(\mathbf{x} - \mathbf{y}_{t+1}) \cdot (\mathbf{x} - \mathbf{v}') \leq \epsilon$ , where  $(\mathbf{v}', \mathbf{s}') \leftarrow \hat{\mathcal{O}}_{\mathcal{K}}(\mathbf{x} - \mathbf{y}_{t+1})$ . Fix  $\mathbf{z} \in CH(\alpha\mathcal{K})$ . It holds that

$$\begin{aligned} \|\mathbf{y}_{t+1} - \mathbf{z}\|^2 &= \|(\mathbf{y}_{t+1} - \mathbf{x}) + (\mathbf{x} - \mathbf{z})\|^2 = \|\mathbf{y}_{t+1} - \mathbf{x}\|^2 + \|\mathbf{x} - \mathbf{z}\|^2 - 2(\mathbf{x} - \mathbf{y}_{t+1}) \cdot (\mathbf{x} - \mathbf{z}) \\ &\geq \|\mathbf{x} - \mathbf{z}\|^2 - 2(\mathbf{x} - \mathbf{y}_{t+1}) \cdot (\mathbf{x} - \mathbf{z}) \\ &\geq \|\mathbf{x} - \mathbf{z}\|^2 - 2(\mathbf{x} - \mathbf{y}_{t+1}) \cdot (\mathbf{x} - \mathbf{v}') \geq \|\mathbf{x} - \mathbf{z}\|^2 - 2\epsilon, \end{aligned}$$

where the second inequality holds since  $\mathbf{v}'$  is the output of the extended approximation oracle with respect to the input  $(\mathbf{x} - \mathbf{y}_{t+1})$ .

For the second part of the lemma, we observe that if  $(\mathbf{x} - \mathbf{y}_{t+1}) \cdot (\mathbf{x} - \mathbf{v}') > \epsilon$ , then

$$\begin{aligned} \|\mathbf{x}' - \mathbf{y}_{t+1}\|^2 &= \|\mathbf{x} - \mathbf{y}_{t+1} + \lambda(\mathbf{v}' - \mathbf{x})\|^2 \\ &= \|\mathbf{x} - \mathbf{y}_{t+1}\|^2 - 2\lambda(\mathbf{x} - \mathbf{y}_{t+1}) \cdot (\mathbf{x} - \mathbf{v}') + \lambda^2\|\mathbf{v}' - \mathbf{x}\|^2 \\ &\leq \|\mathbf{x} - \mathbf{y}_{t+1}\|^2 - 2\lambda\epsilon + 2\lambda^2(\|\mathbf{v}'\|^2 + \|\mathbf{x}\|^2) \\ &\leq \|\mathbf{x} - \mathbf{y}_{t+1}\|^2 - 2\lambda\epsilon + 4\lambda^2(\alpha + 2)^2 R^2, \end{aligned}$$

where the first inequality follows since  $(\mathbf{x} - \mathbf{y}_{t+1}) \cdot (\mathbf{x} - \mathbf{v}') > \epsilon$  and using the triangle inequality with  $(a + b)^2 \leq 2(a^2 + b^2)$ , and the second inequality follows by the assumption on  $\mathbf{x}$  and since  $\mathbf{v}'$  is the output of the extended approximation oracle. Thus, we can see that setting  $\lambda = \frac{\epsilon}{3(\alpha+2)^2 R^2} \in (0, 1]$ , gives the requested result.

Finally, since  $\mathbf{x}$  and  $\mathbf{v}'$  are dominated by  $\mathbf{s}$  and  $\mathbf{s}'$  for any  $\mathbf{f} \in \mathcal{F}$ , respectively, we have that  $\mathbf{x}' = (1 - \lambda)\mathbf{x} + \lambda\mathbf{v}'$  is dominated by  $(1 - \lambda)\mathbf{s} + \lambda\mathbf{s}'$  for any  $\mathbf{f} \in \mathcal{F}$ .  $\square$

Note that Lemma A.1 suggests an iterative algorithm to compute an  $\epsilon$ -approximated projection of  $\mathbf{y}_{t+1}$  in Algorithm 1, that on each iteration reduces the potential  $\|\mathbf{x} - \mathbf{y}_{t+1}\|^2$  by  $\Omega(\epsilon^2)$ , until finding an  $\epsilon$ -approximated projection of  $\mathbf{y}_{t+1}$ ,  $\mathbf{x}_{t+1}$ , which must be found since the potential is non-negative. Moreover, this algorithm finds a point  $\bar{\mathbf{s}}_{t+1} \in CH(\mathcal{K})$ , given explicitly as a convex combination of points in  $\mathcal{K}$  (since  $\lambda \in (0, 1)$ ), such that  $\bar{\mathbf{s}}_{t+1}$  dominates  $\mathbf{x}_{t+1}$  for all vectors in  $\mathcal{F}$ . In particular, sampling  $\mathbf{s}_{t+1}$  from this decomposition guarantees that we play a feasible point in  $\mathcal{K}$ , which in expectation, dominates  $\mathbf{x}_{t+1}$  for all vectors in  $\mathcal{F}$ . The full algorithm, which is closely related to the classical Frank-Wolfe algorithm for convex optimization (a.k.a. the conditional gradient method) [5], is given below, see Algorithm 4<sup>4</sup>.

**Lemma A.2.** *Fix  $\epsilon \in (0, 3(\alpha + 2)^2 R^2]$ ,  $\eta > 0$  and a sequence of loss functions  $\{\mathbf{f}_1, \dots, \mathbf{f}_T\} \subseteq \mathcal{F}$ . Consider the application of Algorithm 1 with learning rate  $\eta$  when applied with respect to the feasible set  $CH(\alpha\mathcal{K})$  and the sequence of losses  $\{\mathbf{f}_1, \dots, \mathbf{f}_T\} \subseteq \mathcal{F}$ , and when we use the algorithm described*

<sup>4</sup>we note that it differs somewhat in presentation than the original algorithm in [13].

above to produce the (randomized) sequence of points  $\{(\mathbf{x}_t, \mathbf{s}_t)\}_{t \in [T]} \subset \mathbb{R}^d \times \mathcal{K}$ . Then, focusing on the case  $\alpha \geq 1$ , it holds that,

$$\begin{aligned} \mathbb{E} \left[ \frac{1}{T} \sum_{t=1}^T \mathbf{s}_t \cdot \mathbf{f}_t \right] - \min_{\mathbf{x} \in \text{CH}(\alpha\mathcal{K})} \frac{1}{T} \sum_{t=1}^T \mathbf{x} \cdot \mathbf{f}_t &= \mathbb{E} \left[ \frac{1}{T} \sum_{t=1}^T \mathbf{s}_t \cdot \mathbf{f}_t \right] - \alpha \min_{\mathbf{x} \in \mathcal{K}} \frac{1}{T} \sum_{t=1}^T \mathbf{x} \cdot \mathbf{f}_t \\ &\leq \frac{\alpha^2 R^2}{T\eta} + \frac{\eta F^2}{2} + \frac{\epsilon}{\eta}. \end{aligned}$$

Moreover, the number of calls to the extended approximation oracle per iteration  $t$  is  $O(\|\mathbf{y}_{t+1} - \mathbf{x}_t\|^2/\epsilon^2) = O(\eta^2 F^2/\epsilon^2)$ , where the  $O(\cdot)$  notation hides polynomial dependencies on  $(1 + \alpha), R$ .

*Proof.* We begin by proving the regret bound. Since each  $\mathbf{x}_{t+1}$  is an approximated projection of  $\mathbf{y}_{t+1}$  in the sense that

$$\forall \mathbf{z} \in \text{CH}(\alpha\mathcal{K}) : \quad \|\mathbf{z} - \mathbf{x}_{t+1}\|^2 \leq \|\mathbf{z} - \mathbf{y}_{t+1}\|^2 + 2\epsilon,$$

it is immediate to see from the proof of Lemma 2.2, that incorporating this approximation error into the regret bound, and bounding  $\|\mathbf{f}_t\| \leq F$  for all  $t$ , results in the regret bound:

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \cdot \mathbf{f}_t - \min_{\mathbf{x} \in \text{CH}(\alpha\mathcal{K})} \frac{1}{T} \sum_{t=1}^T \mathbf{x} \cdot \mathbf{f}_t &= \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \cdot \mathbf{f}_t - \alpha \min_{\mathbf{x} \in \mathcal{K}} \frac{1}{T} \sum_{t=1}^T \mathbf{x} \cdot \mathbf{f}_t \\ &\leq \frac{\alpha^2 R^2}{T\eta} + \frac{\eta F^2}{2} + \frac{\epsilon}{\eta}. \end{aligned}$$

The regret bound now follows by recalling that for all  $t$  and all  $\mathbf{f} \in \mathcal{F}$ :  $\mathbb{E}[\mathbf{s}_t \cdot \mathbf{f}_t] = \bar{\mathbf{s}}_t \cdot \mathbf{f}_t \leq \mathbf{x}_t \cdot \mathbf{f}_t$ , and taking expectation with respect to the randomness in  $\mathbf{s}_t$ .

To bound the number of calls to the approximation oracle per some iteration  $t$ , note that  $\|\mathbf{x}_t - \mathbf{y}_{t+1}\|^2 \leq \eta^2 F^2$ . Thus, if we initialize the projection algorithm, described in Lemma A.1, with the point  $\mathbf{x}_t$ , and we recall that each iteration of the algorithm reduces the potential  $\|\mathbf{x} - \mathbf{y}_{t+1}\|^2$  by  $\Omega(\epsilon^2)$ , where  $\mathbf{x}$  is the current iterate, then we have that at most  $O(\eta^2 F^2/\epsilon^2)$  iterations are required for the algorithm to terminate.  $\square$

The extra term of  $\epsilon/\eta$  in the regret bound is due to fact we compute  $\epsilon$ -approximated projections.

It is clear that setting  $\eta = O(1/\sqrt{T})$  and  $\epsilon = O(1/T)$  in Lemma A.2 guarantees  $O(T^{-1/2})$  expected  $\alpha$ -regret, which is optimal in  $T$ , however requires  $O(T)$  calls to the approximation oracle per iteration. We can also observe that for any constants  $a \in (0, 1), b \geq 1$ , and sufficiently large  $T$ , Lemma A.2 cannot guarantee  $O(T^{-a})$  expected  $\alpha$ -regret using only  $O(\log^b T)$  calls to the approximation oracle per iteration, even on average. For this reason, in this paper we consider a drastically different algorithmic approach to applying the online gradient descent without feasibility methodology.

---

**Algorithm 4** Frank-Wolfe for Approximated (infeasible) Projection onto  $\text{CH}(\alpha\mathcal{K})$

---

- 1: input: point to project  $\mathbf{y} \in \mathbb{R}^d$ , error tolerance  $\epsilon \in (0, 3(\alpha + 2)^2 R^2)$
  - 2: output:  $(\mathbf{x}, \bar{\mathbf{s}}) \in \mathbb{R}^d \times \text{CH}(\mathcal{K})$  such that  $\mathbf{x}$  is an  $\epsilon$ -approximated infeasible projection of  $\mathbf{y}$  dominated by  $\bar{\mathbf{s}}$  for any  $\mathbf{f} \in \mathcal{F}$
  - 3: let  $(\mathbf{x}_1, \bar{\mathbf{s}}_1) \in \mathbb{R}^n \times \text{CH}(\mathcal{K})$  such that  $\mathbf{x}_1$  is dominated by  $\bar{\mathbf{s}}_1$  for any  $\mathbf{f} \in \mathcal{F}$ .
  - 4:  $\lambda \leftarrow \epsilon/(3(\alpha + 2)^2 R^2)$
  - 5: **for**  $i = 1 \dots$  **do**
  - 6:    $(\mathbf{v}_i, \mathbf{s}_i) \leftarrow \hat{\mathcal{O}}_{\mathcal{K}}(\mathbf{x}_i - \mathbf{y})$
  - 7:   **if**  $(\mathbf{x}_i - \mathbf{y}) \cdot (\mathbf{x}_i - \mathbf{v}_i) \leq \epsilon$  **then**
  - 8:     **return**  $(\mathbf{x}_i, \bar{\mathbf{s}}_i)$
  - 9:   **end if**
  - 10:  $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + \lambda(\mathbf{v}_i - \mathbf{x}_i)$
  - 11:  $\bar{\mathbf{s}}_{i+1} \leftarrow \bar{\mathbf{s}}_i + \lambda(\mathbf{s}_i - \bar{\mathbf{s}}_i)$
  - 12: **end for**
-

## B Proofs Omitted from Section 2

### B.1 Proof of Lemma 2.1

*Proof.* For the first item in the lemma note that for  $\alpha \geq 1$  it holds that

$$\begin{aligned}
\mathbf{v} \cdot \mathbf{c} &= \mathcal{O}_{\mathcal{K}}(\mathbf{c}^+) \cdot (\mathbf{c}^+ + \mathbf{c}^-) - \alpha R \bar{\mathbf{c}}^- \cdot (\mathbf{c}^+ + \mathbf{c}^-) \\
&\leq \alpha \min_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot \mathbf{c}^+ + \mathcal{O}_{\mathcal{K}}(\mathbf{c}^+) \cdot \mathbf{c}^- - \alpha R \|\mathbf{c}^-\| \\
&\leq \alpha \min_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot \mathbf{c}^+ - \alpha R \|\mathbf{c}^-\| \\
&\leq \alpha \min_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot \mathbf{c}^+ + \alpha \min_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot \mathbf{c}^- \\
&\leq \alpha \min_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot \mathbf{c} = \min_{\mathbf{x} \in \alpha \mathcal{K}} \mathbf{x} \cdot \mathbf{c},
\end{aligned}$$

Similarly, for  $\alpha > 1$  we have that

$$\begin{aligned}
\mathbf{v} \cdot \mathbf{c} &= \mathcal{O}_{\mathcal{K}}(-\mathbf{c}^-) \cdot (\mathbf{c}^+ + \mathbf{c}^-) - R \bar{\mathbf{c}}^+ \cdot (\mathbf{c}^+ + \mathbf{c}^-) \\
&\leq -\alpha \max_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot (-\mathbf{c}^-) + \mathcal{O}_{\mathcal{K}}(-\mathbf{c}^-) \cdot \mathbf{c}^+ - R \|\mathbf{c}^+\| \\
&\leq \alpha \min_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot \mathbf{c}^- + R \|\mathbf{c}^+\| - R \|\mathbf{c}^+\| \\
&\leq \alpha \min_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot \mathbf{c}^- + \alpha \min_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot \mathbf{c}^+ \\
&\leq \alpha \min_{\mathbf{x} \in \mathcal{K}} \mathbf{x} \cdot \mathbf{c} = \min_{\mathbf{x} \in \alpha \mathcal{K}} \mathbf{x} \cdot \mathbf{c}.
\end{aligned}$$

For the second item, it suffices to observe that for  $\alpha \geq 1$  we have that  $\mathbf{s} \leq \mathbf{v}$  (coordinate-wise) and hence for every  $\mathbf{f} \in \mathcal{F}$  we have that  $\mathbf{s} \cdot \mathbf{f} \leq \mathbf{v} \cdot \mathbf{f}$  (recall that  $\mathcal{F} \subset \mathbb{R}_+^d$ ). Similarly, when  $\alpha < 1$ , we note that  $\mathbf{s} \geq \mathbf{v}$ .

The third item holds trivially.  $\square$

### B.2 Proof of Lemma 2.2

*Proof.* Fix  $\mathbf{x} \in \mathcal{S}$ . Assume that the vectors  $\mathbf{f}_1, \dots, \mathbf{f}_T$  are losses. By the definition of the infeasible projection  $\mathbf{x}_{t+1}$ , for any iteration  $t \geq 1$  it holds that

$$\begin{aligned}
\|\mathbf{x}_{t+1} - \mathbf{x}\|^2 &\leq \|\mathbf{y}_{t+1} - \mathbf{x}\|^2 = \|\mathbf{x}_t - \eta \mathbf{f}_t - \mathbf{x}\|^2 \\
&= \|\mathbf{x}_t - \mathbf{x}\|^2 - 2\eta(\mathbf{x}_t - \mathbf{x}) \cdot \mathbf{f}_t + \eta^2 \|\mathbf{f}_t\|^2
\end{aligned}$$

Rearranging and summing over all iterations we have that

$$\begin{aligned}
\sum_{t=1}^T (\mathbf{x}_t - \mathbf{x}) \cdot \mathbf{f}_t &\leq \frac{1}{2\eta} \sum_{t=1}^T (\|\mathbf{x}_t - \mathbf{x}\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}\|^2) + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{f}_t\|^2 \\
&\leq \frac{1}{2\eta} \|\mathbf{x}_1 - \mathbf{x}\|^2 + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{f}_t\|^2.
\end{aligned}$$

It is immediate to see that the proof of the result in case of payoffs instead of losses (for which the only change is in the update of  $\mathbf{y}_{t+1}$  in Algorithm 1), follows the same lines as the one for losses given above.  $\square$

## C Lemmas and Proofs Omitted from Section 3

### C.1 Proof of Lemma 3.1

*Proof.* To prove the first part of the lemma, suppose there exists some iteration during which the Ellipsoid method declares Problem (4) feasible, and let  $\mathbf{w}$  be the corresponding iterate and let  $(\mathbf{v}, \mathbf{s})$  be the output of the extended approximation oracle on that iteration. Clearly it holds that

$$\epsilon \leq (\mathbf{x} - \mathbf{v}) \cdot \mathbf{w} = \mathbf{x} \cdot \mathbf{w} + \mathbf{v} \cdot (-\mathbf{w}) \leq \mathbf{x} \cdot \mathbf{w} + \min_{\mathbf{z} \in \alpha \mathcal{K}} \mathbf{z} \cdot (-\mathbf{w}) = \min_{\mathbf{z} \in \alpha \mathcal{K}} (\mathbf{x} - \mathbf{z}) \cdot \mathbf{w},$$

where the first inequality follows from the fact that the Ellipsoid method declared Problem (4) feasible, and the second inequality follows from the definition of the extended approximation oracle. Since the Ellipsoid method declared Problem (4) feasible, it also follows that  $\|\mathbf{w}\| \leq 1$  and hence  $\mathbf{w}$  is indeed a feasible solution to Problem (4).

Consider now the case that all  $N$  iterations are executed without declaring Problem (4) feasible and let  $\mathbf{v}_1, \dots, \mathbf{v}_N$  be as defined in the lemma. We would like to show that this implies that

$$\forall \text{ unit vector } \mathbf{w} : \min_{i \in \{1, \dots, N'\}} (\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{w} \leq 3\epsilon. \quad (6)$$

Then, the second part of the lemma follows from applying the next lemma, Lemma C.1, which shows that (6) implies that the point  $\mathbf{p}$  defined in the lemma indeed satisfies  $\|\mathbf{p} - \mathbf{x}\| \leq 3\epsilon$ , as required.

Towards proving (6), suppose that there exists a unit vector  $\mathbf{h} \in \mathbb{R}^d$  such that for all  $i \in \{1, \dots, N\}$ ,  $(\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{h} > 3\epsilon$ . It follows that  $\forall i \in \{1, \dots, N\} : (\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{h}/2 > 3\epsilon/2$ . It follows from a simple application of the Cauchy-Swartz inequality and the observation that  $\|\mathbf{x} - \mathbf{v}_i\| \leq \|\mathbf{x}\| + \|\mathbf{v}_i\| \leq \|\mathbf{x}\| + (\alpha + 2)R$ , that denoting  $r := \frac{\epsilon}{2(\alpha+2)R + \|\mathbf{x}\|}$ , we have that

$$\forall \mathbf{h}' \in \mathcal{B}(\mathbf{h}/2, r) : \min_{i \in [N]} (\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{h}' > \epsilon. \quad (7)$$

Note that on one hand, by the above and our assumption on  $\epsilon$ , every point in  $\mathcal{B}(\mathbf{h}/2, r)$  satisfies the stopping criteria of the Ellipsoid method described in the lemma. On the other-hand, on every iteration in which the current iterate  $\mathbf{w}$  is not declared feasible, it follows that the separating hyperplane fed to the Ellipsoid method indeed separates  $\mathbf{w}$  from  $\mathcal{B}(\mathbf{h}/2, r)$ . To see why this is true, we consider the two possible options for the separating hyperplane. If the hyperplane is  $\mathbf{v}_i - \mathbf{x}$ , where  $\mathbf{v}_i$  is the output of the extended approximation oracle on that iteration, then we have that

$$\forall \mathbf{h}' \in \mathcal{B}(\mathbf{h}/2, r) : (\mathbf{w} - \mathbf{h}') \cdot (\mathbf{v}_i - \mathbf{x}) = (\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{h}' - (\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{w} > \epsilon - \epsilon = 0,$$

where the first inequality follows from Eq. (7) and the fact that  $(\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{w} < \epsilon$  on this iteration. If the hyperplane used was  $\mathbf{w}$ , which guarantees that on that iteration  $\|\mathbf{w}\| > 1$ , then we have that

$$\forall \mathbf{h}' \in \mathcal{B}(\mathbf{h}/2, r) : (\mathbf{w} - \mathbf{h}') \cdot \mathbf{w} = \|\mathbf{w}\| - \mathbf{w} \cdot \mathbf{h}' > 1 - 1 = 0,$$

where the last inequality follows since by our assumption on  $\epsilon$ , it holds that  $\mathcal{B}(\mathbf{h}/2, r) \subset \mathcal{B}(0, 1)$ . Thus, we can conclude that if the number of Ellipsoid method iterations satisfies  $N \geq cd^2 \ln \left( \frac{(\alpha+1)R + \|\mathbf{x}\|}{\epsilon} \right)$  for an appropriate universal constant  $c > 0$ , and all  $N$  iterations were completed without declaring feasibility, it follows that no such unit vector  $\mathbf{h}$  can exist, which means Eq. (6) holds, and the result follows.  $\square$

**Lemma C.1.** Fix  $\mathbf{x} \in \mathbb{R}^d$ , vectors  $\mathbf{v}_1, \dots, \mathbf{v}_N \in \mathbb{R}^d$  and  $\epsilon > 0$ . If for any unit vector  $\mathbf{w}$  it holds that  $\min_{i \in \{1, \dots, N\}} (\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{w} \leq \epsilon$ , then it follows that the point  $\mathbf{p} = \sum_{i=1}^N a_i \mathbf{v}_i$ , where  $(a_1, \dots, a_N)$  is an optimal solution to Problem (5), satisfies  $\|\mathbf{p} - \mathbf{x}\| \leq \epsilon$ .

*Proof.* First we show that the following holds:

$$\begin{aligned} \forall i, j \text{ s.t. } a_i > 0, a_j > 0 : & (\mathbf{p} - \mathbf{x}) \cdot \mathbf{v}_i = (\mathbf{p} - \mathbf{x}) \cdot \mathbf{v}_j, \\ \forall i, j \text{ s.t. } a_i > 0, a_j = 0 : & (\mathbf{p} - \mathbf{x}) \cdot \mathbf{v}_i \leq (\mathbf{p} - \mathbf{x}) \cdot \mathbf{v}_j. \end{aligned} \quad (8)$$

To see why this is true, fix some  $i, j$  such that  $a_i > 0$  and consider the point  $\mathbf{p}' = \mathbf{p} + \delta(\mathbf{v}_j - \mathbf{v}_i)$  such that  $0 < \delta \leq a_i$ . Clearly  $\mathbf{p}'$  lies in the convex hull of  $\{\mathbf{v}_1, \dots, \mathbf{v}_N\}$  and hence is a feasible solution to Problem (5). It holds that

$$\frac{1}{2} \|\mathbf{p}' - \mathbf{x}\|^2 = \frac{1}{2} \|\mathbf{p} - \mathbf{x}\|^2 + \delta(\mathbf{v}_j - \mathbf{v}_i) \cdot (\mathbf{p} - \mathbf{x}) + \frac{\delta^2}{2} \|\mathbf{v}_i - \mathbf{v}_j\|^2. \quad (9)$$

Thus, we can see that if (8) does not hold, then without loss of generality we can always choose  $i, j$  such that  $a_i > 0$  and  $(\mathbf{p} - \mathbf{x}) \cdot \mathbf{v}_i > (\mathbf{p} - \mathbf{x}) \cdot \mathbf{v}_j$ , and thus as can be seen from Eq. (9), choosing  $\delta$  to be sufficiently small it follows that  $\|\mathbf{p}' - \mathbf{x}\|^2 < \|\mathbf{p} - \mathbf{x}\|^2$ , contradicting the optimality of  $\mathbf{p}$ .

Denoting by  $\mathbf{u}$  the unit vector in the direction of  $\mathbf{x} - \mathbf{p}$ , we can rewrite Eq. (8) as follows:

$$\begin{aligned} \forall i, j \text{ s.t. } a_i > 0, a_j > 0 : & (\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{u} = (\mathbf{x} - \mathbf{v}_j) \cdot \mathbf{u}, \\ \forall i, j \text{ s.t. } a_i > 0, a_j = 0 : & (\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{u} \leq (\mathbf{x} - \mathbf{v}_j) \cdot \mathbf{u}. \end{aligned} \quad (10)$$

Using our assumption, we in particular have that  $\min_{i \in [N]} (\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{u} \leq \epsilon$ , and using Eq. (10) we have that

$$\|\mathbf{p} - \mathbf{x}\| = (\mathbf{x} - \mathbf{p}) \cdot \mathbf{u} = \sum_{i=1}^N a_i (\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{u} = \min_{i \in [N]} (\mathbf{x} - \mathbf{v}_i) \cdot \mathbf{u} \leq \epsilon,$$

where the last equality is a consequence of Eq. (10) and the fact that  $(a_1, \dots, a_N)$  is a distribution. Thus the lemma follows.  $\square$

## C.2 Proof of Lemma 3.2

*Proof.* Note that the second item in the lemma is a straightforward guarantee of Lemma 3.1.

To prove the first item, suppose that the algorithm terminates after the **for** loop was entered  $k$  times, and let  $\tilde{\mathbf{y}}_1, \dots, \tilde{\mathbf{y}}_k$  denote the values of  $\tilde{\mathbf{y}}$  throughout the run of the algorithm, where  $\tilde{\mathbf{y}}_i$  is the value of  $\tilde{\mathbf{y}}$  at the beginning of the  $i$ th iteration of the **for** loop. Note that since  $\text{CH}(\alpha\mathcal{K}) \subseteq \mathcal{B}(0, \alpha R)$  and  $\tilde{\mathbf{y}}_1$  is the projection of  $\mathbf{y}$  onto  $\mathcal{B}(0, \alpha R)$ , we have that  $\forall \mathbf{z} \in \text{CH}(\alpha\mathcal{K}) : \|\tilde{\mathbf{y}}_1 - \mathbf{z}\|^2 \leq \|\mathbf{y} - \mathbf{z}\|^2$ .

We are now going to show that for any  $i \geq 1$  it holds that

$$\forall \mathbf{z} \in \text{CH}(\alpha\mathcal{K}) : \|\tilde{\mathbf{y}}_{i+1} - \mathbf{z}\|^2 \leq \|\tilde{\mathbf{y}}_i - \mathbf{z}\|^2, \quad (11)$$

which clearly yields item 1 in the lemma.

To prove that Eq. (11) holds throughout the run of the algorithm, consider an iteration  $i$  of the **for** loop during which, the SEPARATION-OR-DECOMPOSITION procedure returns a separating hyperplane  $\mathbf{w}$ . It holds that

$$\begin{aligned} \forall \mathbf{z} \in \text{CH}(\alpha\mathcal{K}) : \|\tilde{\mathbf{y}}_i - \mathbf{z}\|^2 &= \|\tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_{i+1} + \tilde{\mathbf{y}}_{i+1} - \mathbf{z}\|^2 \\ &= \|\tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_{i+1}\|^2 + \|\tilde{\mathbf{y}}_{i+1} - \mathbf{z}\|^2 + 2(\tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_{i+1}) \cdot (\tilde{\mathbf{y}}_{i+1} - \mathbf{z}) \\ &\geq \|\tilde{\mathbf{y}}_{i+1} - \mathbf{z}\|^2 + 2(\tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_{i+1}) \cdot (\tilde{\mathbf{y}}_{i+1} - \mathbf{z}) \\ &= \|\tilde{\mathbf{y}}_{i+1} - \mathbf{z}\|^2 + 2\epsilon \mathbf{w} \cdot [(\tilde{\mathbf{y}}_i - \mathbf{z}) - \epsilon \mathbf{w}] \\ &= \|\tilde{\mathbf{y}}_{i+1} - \mathbf{z}\|^2 + 2\epsilon(\tilde{\mathbf{y}}_i - \mathbf{z}) \cdot \mathbf{w} - 2\epsilon^2 \|\mathbf{w}\|^2 \\ &\geq \|\tilde{\mathbf{y}}_{i+1} - \mathbf{z}\|^2, \end{aligned}$$

where the third equality follows from the update rule of  $\tilde{\mathbf{y}}$  in the algorithm, and the last inequality is a direct consequence of the guarantees of Lemma 3.1. Thus, Eq. (11) indeed holds for all  $i \geq 1$ , which gives the first item listed in the lemma.

We now turn to upper bound the number of iterations performed by the algorithm. Consider again an iteration  $i$  of the loop during which the SEPARATION-OR-DECOMPOSITION procedure returns a separating hyperplane  $\mathbf{w}$ . We are going to show that

$$\text{dist}^2(\tilde{\mathbf{y}}_{i+1}, \text{CH}(\alpha\mathcal{K})) \leq \text{dist}^2(\tilde{\mathbf{y}}_i, \text{CH}(\alpha\mathcal{K})) - \epsilon^2,$$

which, together with the fact that  $\text{dist}^2(\tilde{\mathbf{y}}_1, \text{CH}(\alpha\mathcal{K})) \leq 2\alpha^2 R^2$ , gives the desired upper bound on the number of iterations.

Denote  $\mathbf{x}_i = \arg \min_{\mathbf{x} \in \text{CH}(\alpha\mathcal{K})} \|\mathbf{x} - \tilde{\mathbf{y}}_i\|$  and  $\mathbf{x}_{i+1} = \arg \min_{\mathbf{x} \in \text{CH}(\alpha\mathcal{K})} \|\mathbf{x} - \tilde{\mathbf{y}}_{i+1}\|$ . It holds that

$$\begin{aligned} \text{dist}^2(\tilde{\mathbf{y}}_{i+1}, \text{CH}(\alpha\mathcal{K})) &= \|\mathbf{x}_{i+1} - \tilde{\mathbf{y}}_{i+1}\|^2 \leq \|\mathbf{x}_i - \tilde{\mathbf{y}}_{i+1}\|^2 = \|\mathbf{x}_i - \tilde{\mathbf{y}}_i + \epsilon \mathbf{w}\|^2 \\ &= \text{dist}^2(\tilde{\mathbf{y}}_i, \text{CH}(\alpha\mathcal{K})) + \epsilon^2 \|\mathbf{w}\|^2 - 2\epsilon(\tilde{\mathbf{y}}_i - \mathbf{x}_i) \cdot \mathbf{w} \\ &\leq \text{dist}^2(\tilde{\mathbf{y}}_i, \text{CH}(\alpha\mathcal{K})) - \epsilon^2, \end{aligned}$$

where the inequality is a direct consequence of the guarantees of Lemma 3.1. Thus, we obtain both the desired bound on the number of iterations and the bound on the distance of the final point  $\tilde{\mathbf{y}}$  from  $\text{CH}(\alpha\mathcal{K})$ .

Finally, we turn to upper bound to overall number of queries to the approximation oracle. Using the bound in Lemma 3.1, we have that the number of calls to the oracle on the  $i$ th iteration of the loop is

upper bounded by  $O\left(d^2 \ln\left(\frac{(\alpha+1)R + \|\tilde{\mathbf{y}}_i\|}{\epsilon}\right)\right)$ . As we have shown, the values  $\text{dist}(\tilde{\mathbf{y}}_i, \text{CH}(\alpha\mathcal{K}))$  are monotonically decreasing with  $i$  and hence we can upper bound

$$\|\tilde{\mathbf{y}}_i\| \leq \max_{\mathbf{x} \in \text{CH}(\alpha\mathcal{K})} \|\mathbf{x}\| + \text{dist}(\tilde{\mathbf{y}}_i, \text{CH}(\alpha\mathcal{K})) \leq \alpha R + \text{dist}(\tilde{\mathbf{y}}_1, \text{CH}(\alpha\mathcal{K})) \leq \alpha R + \sqrt{2}\alpha R,$$

where the last inequality holds since  $\tilde{\mathbf{y}}_1$  is the projection of  $\mathbf{y}$  onto the ball  $\mathcal{B}(0, \alpha R)$ . Thus, the overall number of queries to the approximation oracle after  $k$  iterations is upper bounded by  $O\left(kd^2 \ln\left(\frac{(\alpha+1)R}{\epsilon}\right)\right)$ .  $\square$

## D Algorithms and Proofs Omitted from Section 4

### D.1 Proof of Theorem 4.1

*Proof.* For the proof we focus on the case  $\alpha \geq 1$  since the proof for the complementary follows from the same derivations up to changes in the obvious places. To prove the regret bound, we simply apply Lemma 2.2 with respect to the sequence of points  $\{\tilde{\mathbf{y}}_t\}_{t=1}^T$  and the feasible set  $\text{CH}(\alpha\mathcal{K})$  and plugin the guarantee of Lemma 3.2, which gives

$$\begin{aligned} \sum_{t=1}^T \tilde{\mathbf{y}}_t \cdot \mathbf{f}_t - \min_{\mathbf{x} \in \alpha\mathcal{K}} \sum_{t=1}^T \mathbf{x} \cdot \mathbf{f}_t &= \sum_{t=1}^T \tilde{\mathbf{y}}_t \cdot \mathbf{f}_t - \alpha \cdot \min_{\mathbf{x} \in \mathcal{K}} \sum_{t=1}^T \mathbf{x} \cdot \mathbf{f}_t \\ &\leq \frac{\alpha^2 R^2}{\eta} + T \frac{\eta F^2}{2}, \end{aligned}$$

where we have used the fact that  $\|\tilde{\mathbf{y}}_1\| \leq \alpha R$  and  $\|\mathbf{f}_t\| \leq F$  for all  $t \in [T]$ . For every iteration  $t \geq 1$ , let us denote  $\mathbf{p}_{t+1} = \sum_{i=1}^N a_i \mathbf{v}_i$ ,  $\bar{\mathbf{s}}_t = \sum_{i=1}^N a_i \mathbf{s}_i$ , where  $(a_1, \dots, a_N)$ ,  $\{(\mathbf{v}_1, \mathbf{s}_1), \dots, (\mathbf{v}_N, \mathbf{s}_N)\}$  are the outputs of the call to Algorithm 2 on iteration  $t$ , and for  $t = 1$  we denote  $\mathbf{p}_1 = \tilde{\mathbf{y}}_1$  and  $\bar{\mathbf{s}}_1 = \mathbf{s}_1$ . By the guarantee of Lemma 3.2, we have that

$$\sum_{t=1}^T \mathbf{p}_t \cdot \mathbf{f}_t - \alpha \cdot \min_{\mathbf{x} \in \mathcal{K}} \sum_{t=1}^T \mathbf{x} \cdot \mathbf{f}_t \leq \frac{\alpha^2 R^2}{\eta} + T \frac{\eta F^2}{2} + 3T\epsilon F,$$

where the inequality holds since for all  $t \geq 1$ :  $|(\mathbf{p}_t - \tilde{\mathbf{y}}_t) \cdot \mathbf{f}_t| \leq \|\mathbf{p}_t - \tilde{\mathbf{y}}_t\| \cdot \|\mathbf{f}_t\| \leq 3\epsilon F$ . The regret bound now follows since for any iteration  $t$ ,  $\bar{\mathbf{s}}_t$  dominates  $\mathbf{p}_t$  for any vector  $\mathbf{f} \in \mathcal{F}$ , and since  $\mathbb{E}[\mathbf{s}_t] = \bar{\mathbf{s}}_t$ .

We now turn to upper bound the overall number of queries to the approximation oracle of  $\mathcal{K}$ . Let  $k_t$  be the number of iterations it took Algorithm 2 to terminate, when invoked on iteration  $t$  of Algorithm 3. Note that, by Lemma 3.2, we have that  $K(\eta, \epsilon) = O\left(\frac{1}{T} \sum_{t=1}^{T-1} k_t d^2 \ln\left(\frac{(\alpha+1)R}{\epsilon}\right)\right)$ . By Lemma 3.2, it follows that on any iteration  $t$ ,

$$\begin{aligned} \text{dist}^2(\tilde{\mathbf{y}}_{t+1}, \text{CH}(\alpha\mathcal{K})) &\leq \text{dist}^2(\mathbf{y}_{t+1}, \text{CH}(\alpha\mathcal{K})) - (k_t - 1)\epsilon^2 \\ &= \text{dist}^2(\tilde{\mathbf{y}}_t - \eta \mathbf{f}_t, \text{CH}(\alpha\mathcal{K})) - (k_t - 1)\epsilon^2 \\ &\leq (\text{dist}(\tilde{\mathbf{y}}_t, \text{CH}(\alpha\mathcal{K})) + \eta F)^2 - (k_t - 1)\epsilon^2 \\ &= \text{dist}^2(\tilde{\mathbf{y}}_t, \text{CH}(\alpha\mathcal{K})) + 2\eta F \text{dist}(\tilde{\mathbf{y}}_t, \text{CH}(\alpha\mathcal{K})) + \eta^2 F^2 - k_t \epsilon^2 + \epsilon^2. \end{aligned}$$

Rearranging, summing over all  $T$  iterations, and recalling that for all  $t$ ,  $\text{dist}(\tilde{\mathbf{y}}_t, \text{CH}(\alpha\mathcal{K})) \leq \sqrt{2}\alpha R$ , we have that

$$\begin{aligned} \sum_{t=1}^{T-1} k_t &\leq \frac{1}{\epsilon^2} \left( \text{dist}^2(\tilde{\mathbf{y}}_1, \text{CH}(\alpha\mathcal{K})) - \text{dist}^2(\tilde{\mathbf{y}}_T, \text{CH}(\alpha\mathcal{K})) + (T-1) \left( 2\sqrt{2}\eta\alpha RF + \eta^2 F^2 + \epsilon^2 \right) \right) \\ &\leq (T-1) \left( 1 + \frac{2\sqrt{2}\eta\alpha RF + \eta^2 F^2}{\epsilon^2} \right). \end{aligned}$$

$\square$

## D.2 Algorithm for the bandit setting and proof of Theorem 4.2

---

### Algorithm 5 Bandit Algorithm

---

```

1: input: learning rate  $\eta > 0$ , projection error parameter  $\epsilon > 0$ ,  $\{\mathbf{q}_1, \dots, \mathbf{q}_d\}$  - a  $\beta$ -BS( $\mathcal{K}$ ) for some
    $\beta > 0$ , exploration parameter  $\gamma \in (0, 1)$ 
2: instantiate Algorithm 3 with parameters  $(\eta, \epsilon)$ 
3: for  $t = 1 \dots T$  do
4:   receive  $(\mathbf{s}_t, \tilde{\mathbf{y}}_t) \in \mathcal{K} \times \mathcal{B}(0, \alpha R)$  from Algorithm 3
5:    $b_t \leftarrow \begin{cases} \text{EXPLORE} & \text{with prob. } \gamma \\ \text{EXPLOIT} & \text{with prob. } 1 - \gamma \end{cases}$ 
6:   if  $b_t = \text{EXPLORE}$  then
7:     sample  $i_t \in [d]$  uniformly at random
8:     play  $\hat{\mathbf{s}}_t = \mathbf{q}_{i_t}$ 
9:     receive loss/payoff  $\ell_t = \mathbf{q}_{i_t} \cdot \mathbf{f}_t$ 
10:    set  $\hat{\mathbf{f}}_t \leftarrow \frac{d\ell_t}{\gamma} \mathbf{Q}^{-1} \mathbf{q}_{i_t}$  {recall  $\mathbf{Q} = \sum_{i=1}^d \mathbf{q}_i \mathbf{q}_i^\top$ }
11:   else
12:     play  $\hat{\mathbf{s}}_t = \mathbf{s}_t$ 
13:     receive loss/payoff  $\ell_t = \mathbf{s}_t \cdot \mathbf{f}_t$ 
14:     set  $\hat{\mathbf{f}}_t \leftarrow \mathbf{0}$ 
15:   end if
16:   feed  $\hat{\mathbf{f}}_t$  to Algorithm 3 as the loss/payoff vector for round  $t$ 
17: end for

```

---

*proof of Theorem 4.2.* The proof is very similar to that of Theorem 4.1 and we focus on the modifications of it required to prove Theorem 4.2. Again, we focus on the case  $\alpha \geq 1$  since the complementary case follows the same lines with the obvious modifications. Let  $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathcal{K}} \sum_{t=1}^T \mathbf{x} \cdot \mathbf{f}_t$ . Applying Lemma 2.2 with respect to the sequence of points  $\{\tilde{\mathbf{y}}_t\}_{t=1}^T$  and the sequence of losses  $\{\hat{\mathbf{f}}_t\}_{t=1}^T$ , we have that

$$\sum_{t=1}^T \tilde{\mathbf{y}}_t \cdot \hat{\mathbf{f}}_t - \alpha \cdot \sum_{t=1}^T \mathbf{x}^* \cdot \hat{\mathbf{f}}_t \leq \frac{\alpha^2 R^2}{\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\hat{\mathbf{f}}_t\|^2.$$

Taking expectation with respect to the random variables  $b_1, i_1, \dots, b_T, i_T$  and noting that for all  $t \in [T]$ , both  $\mathbf{x}^*$  and  $\tilde{\mathbf{y}}_t$  are independent of the randomness in  $\hat{\mathbf{f}}_t$ , we have that

$$\mathbb{E}_{\{(b_t, i_t)\}_{t=1}^T} \left[ \sum_{t=1}^T \tilde{\mathbf{y}}_t \cdot \hat{\mathbf{f}}_t \right] - \alpha \cdot \sum_{t=1}^T \mathbf{x}^* \cdot \hat{\mathbf{f}}_t \leq \frac{\alpha^2 R^2}{\eta} + T \frac{\eta d^2 C^2 \beta^2}{2\gamma},$$

where we have used the observations that

$$\begin{aligned} \mathbb{E}_{b_t, i_t} [\hat{\mathbf{f}}_t] &= \gamma \sum_{i=1}^d \frac{1}{d} \cdot \frac{d \mathbf{q}_i^\top \mathbf{f}_t}{\gamma} \mathbf{Q}^{-1} \mathbf{q}_i = \sum_{i=1}^d \mathbf{Q}^{-1} \mathbf{q}_i \mathbf{q}_i^\top \mathbf{f}_t = \mathbf{Q}^{-1} \mathbf{Q} \mathbf{f}_t = \mathbf{f}_t, \\ \mathbb{E}_{b_t} [\|\hat{\mathbf{f}}_t\|^2] &= \gamma \frac{d^2}{\gamma^2} \ell_t^2 \|\mathbf{Q}^{-1} \mathbf{q}_{i_t}\|^2 + (1 - \gamma) 0 \leq \frac{(dC\beta)^2}{\gamma}. \end{aligned}$$

As in the proof of Theorem 4.1, for every iteration  $t \geq 1$ , let us denote  $\mathbf{p}_{t+1} = \sum_{i=1}^N a_i \mathbf{v}_i$ ,  $\bar{\mathbf{s}}_t = \sum_{i=1}^N a_i \mathbf{s}_i$ , where  $(a_1, \dots, a_N)$ ,  $\{(\mathbf{v}_1, \mathbf{s}_1), \dots, (\mathbf{v}_N, \mathbf{s}_N)\}$  are the outputs of the call to Algorithm 2 on that iteration. Also define  $\mathbf{p}_1 = \tilde{\mathbf{y}}_1$ . Again, by the guarantee of Lemma 3.2, we have that

$$\mathbb{E}_{\{(b_t, i_t)\}_{t=1}^T} \left[ \sum_{t=1}^T \mathbf{p}_t \cdot \mathbf{f}_t \right] - \alpha \cdot \sum_{t=1}^T \mathbf{x}^* \cdot \mathbf{f}_t \leq \frac{\alpha^2 R^2}{\eta} + T \frac{\eta d^2 C^2 \beta^2}{2\gamma} + 3T\epsilon F.$$

Since  $\mathbf{p}_t$  is dominated by  $\bar{\mathbf{s}}_t = \mathbb{E}[\mathbf{s}_t]$  for all  $t \in [T]$ , we have that

$$\mathbb{E}_{\{(b_t, i_t, \mathbf{s}_t)\}_{t=1}^T} \left[ \sum_{t=1}^T \mathbf{s}_t \cdot \mathbf{f}_t \right] - \alpha \cdot \sum_{t=1}^T \mathbf{x}^* \cdot \mathbf{f}_t \leq \frac{\alpha^2 R^2}{\eta} + T \frac{\eta d^2 C^2 \beta^2}{2\gamma} + 3T\epsilon F.$$



Finally, since

$$\forall t \in [T] : \quad \mathbb{E}_{b_t}[\hat{\mathbf{s}}_t \cdot \mathbf{f}_t] = (1 - \gamma)\mathbf{s}_t \cdot \mathbf{f}_t + \gamma \mathbf{q}_{i_t} \cdot \mathbf{f}_t \begin{cases} \leq \mathbf{s}_t \cdot \mathbf{f}_t + \gamma C & \text{if } \alpha \geq 1 \\ \geq \mathbf{s}_t \cdot \mathbf{f}_t - \gamma C & \text{if } \alpha < 1 \end{cases},$$

we have that

$$\mathbb{E} \left[ \sum_{t=1}^T \hat{\mathbf{s}}_t \cdot \mathbf{f}_t \right] - \alpha \cdot \sum_{t=1}^T \mathbf{x}^* \cdot \mathbf{f}_t \leq \frac{\alpha^2 R^2}{\eta} + T \frac{\eta d^2 C^2 \beta^2}{2\gamma} + 3T\epsilon F + T\gamma C,$$

as required.

We now turn to upper bound the overall number of queries to the approximation oracle of  $\mathcal{K}$ . Note that we require to compute a new approximated projection only after rounds for which it holds that  $b_t = \text{EXPLORE}$ , since otherwise it holds that  $\hat{\mathbf{f}}_t = \mathbf{0}$ , and there is no update to the iterates maintained by Algorithm 3. For any  $t \in [T]$  we define the indicator variable:

$$I_t \leftarrow \begin{cases} 1 & \text{if } b_t = \text{EXPLORE}; \\ 0 & \text{if } b_t = \text{EXPLOIT}. \end{cases}$$

Define  $\hat{F} := \frac{dC\beta}{\gamma}$ , and observe that for all  $t \in [T]$  it holds that

$$\|\hat{\mathbf{f}}_t\| \leq \left\| \frac{d\mathbf{Q}^{-1}\mathbf{q}_{i_t}\ell_t}{\gamma} \right\| \leq \frac{d}{\gamma} |\mathbf{q}_{i_t} \cdot \mathbf{f}_t| \cdot \|\mathbf{Q}^{-1}\mathbf{q}_{i_t}\| \leq \frac{d}{\gamma} C\beta = \hat{F}.$$

Now, we continue to bound the number of calls to Algorithm 2, very similarly to the analysis in the proof of Theorem 4.1.

Let  $k_t$  be the number of iterations it took Algorithm 2 to terminate when invoked on iteration  $t$  of Algorithm 3 (w.l.o.g. this happens when Algorithm 5 sends the feedback  $\hat{\mathbf{f}}_t$  to Algorithm 3), and note that  $\mathbb{E}[K(\eta, \epsilon, \gamma)] = \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^{T-1} k_t \right] \cdot O \left( d^2 \ln \left( \frac{(\alpha+1)R}{\epsilon} \right) \right)$ . Note that for all  $t \geq 1$ ,  $\mathbf{y}_{t+1} = \tilde{\mathbf{y}}_t - I_t \eta \hat{\mathbf{f}}_t$ . Thus, by Lemma 3.2, it follows that on any iteration  $t$ ,

$$\begin{aligned} \text{dist}^2(\tilde{\mathbf{y}}_{t+1}, \text{CH}(\alpha\mathcal{K})) &\leq \text{dist}^2(\mathbf{y}_{t+1}, \text{CH}(\alpha\mathcal{K})) - (k_t - 1)\epsilon^2 \\ &= \text{dist}^2(\tilde{\mathbf{y}}_t - I_t \eta \hat{\mathbf{f}}_t, \text{CH}(\alpha\mathcal{K})) - (k_t - 1)\epsilon^2 \\ &\leq (\text{dist}^2(\tilde{\mathbf{y}}_t, \text{CH}(\alpha\mathcal{K})) + I_t \eta \hat{F})^2 - (k_t - 1)\epsilon^2 \\ &= \text{dist}^2(\tilde{\mathbf{y}}_t, \text{CH}(\alpha\mathcal{K})) + 2I_t \eta \hat{F} \text{dist}(\tilde{\mathbf{y}}_t, \text{CH}(\alpha\mathcal{K})) + I_t \eta^2 \hat{F}^2 - k_t \epsilon^2 + \epsilon^2. \end{aligned}$$

Rearranging, summing over all iterations  $1 \dots T-1$ , and recalling that for all  $t$ ,  $\text{dist}(\tilde{\mathbf{y}}_{t-1}, \text{CH}(\alpha\mathcal{K})) \leq \sqrt{2}\alpha R$ , we have that

$$\sum_{t=1}^{T-1} k_t \leq \frac{1}{\epsilon^2} \left( \text{dist}^2(\tilde{\mathbf{y}}_1, \text{CH}(\alpha\mathcal{K})) - \text{dist}^2(\tilde{\mathbf{y}}_T, \text{CH}(\alpha\mathcal{K})) + 2\sqrt{2} \sum_{t=1}^{T-1} I_t \eta \alpha \hat{F} R + \sum_{t=1}^{T-1} I_t \eta^2 \hat{F}^2 + (T-1)\epsilon^2 \right).$$

Taking expectation with respect to the random variables  $I_1, \dots, I_{T-1}$  we have that

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^{T-1} k_t \right] &\leq (T-1) \left( 1 + \frac{2\sqrt{2}\gamma\eta\alpha\hat{F}R + \gamma\eta^2\hat{F}^2}{\epsilon^2} \right) \\ &= (T-1) \left( 1 + \frac{2\sqrt{2}\eta\alpha\beta dCR + (\eta dC\beta)^2/\gamma}{\epsilon^2} \right), \end{aligned}$$

as required. □