

# Efficient Projection-Free Online Convex Optimization with Membership Oracle

**Zakaria Mhammedi**

*Massachusetts Institute of Technology*

MHAMMEDI@MIT.EDU

**Editors:** Po-Ling Loh and Maxim Raginsky

## Abstract

In constrained convex optimization, existing interior point methods do not scale well with the dimension of the ambient space. Alternative approaches such as Projected Gradient Descent only provide a computational benefit for simple convex sets where Euclidean projections can be performed efficiently, such as Euclidean balls. For other more complex sets, the cost of the projections can be too high. To circumvent these issues, alternative methods based on the famous Frank-Wolfe algorithm have been studied and widely used. Such methods use a Linear Optimization Oracle at each iteration instead of Euclidean projections; the former can often be performed efficiently. Such methods have also been extended to the online and stochastic optimization settings. However, the Frank-Wolfe algorithm and its variants do not achieve the optimal performance, in terms of regret or rate, for general convex sets. What is more, the Linear Optimization Oracle they use can still be computationally expensive in some cases. In this paper, we move away from Frank-Wolfe style algorithms and present a new reduction that turns any algorithm  $A$  over a Euclidean ball (where projections are cheap) to an algorithm over a general convex constraint set  $\mathcal{C}$  contained within the ball, without sacrificing the performance of the original algorithm  $A$  by much. Our reduction requires  $O(T \ln T)$  calls to a Membership Oracle on  $\mathcal{C}$  after  $T$  rounds, and no linear optimization on  $\mathcal{C}$  is needed. Using this reduction, we recover optimal regret bounds [resp. rates], in terms of the number of iterations, in online [resp. stochastic] convex optimization. Our guarantees are also useful in the offline convex optimization setting when the dimension of the ambient space is large.

## 1. Introduction

In this paper, we are interested in designing efficient algorithms for constrained convex optimization in the online, offline, and stochastic settings. Popular algorithms for optimizing a convex objective  $f$  defined on a bounded convex set  $\mathcal{C} \subset \mathbb{R}^d$  include, for example, the ellipsoid or cutting plane methods (Grötschel et al., 1993, 2012; Bubeck, 2015). Though such algorithms enjoy linear convergence rates, where the optimality gap decreases exponentially fast with the number of iterations, their per-iteration computational complexity depends super-linearly in the dimension  $d$ . In particular, if  $\text{Cost}(\mathcal{M}_{\mathcal{C}})$  is the computational cost of testing if some point  $\mathbf{x} \in \mathbb{R}^d$  belongs to  $\mathcal{C}$ , then state-of-the-art cutting plane-based methods have a computational complexity of order  $O(d^3 + d^2 \text{Cost}(\mathcal{M}_{\mathcal{C}})) \log(1/\epsilon)$  to find an  $\epsilon$ -suboptimal point (Lee et al., 2015, 2018). Thus, when the dimension  $d$  is large, such methods can become impractical.

For some convex sets like the Euclidean ball, alternative methods, such as the projected Gradient Descent (GD) algorithm, dispense of the “expensive” dimension dependence in their computational complexity at the cost of a worse dependence in  $1/\epsilon$  (e.g.  $1/\epsilon^2$  instead of  $\log(1/\epsilon)$ ) in their convergence rate. Despite this cost, such methods are still favorable in practice when the dimension  $d$  is large. We refer the interested reader to Bubeck (2015) for a comprehensive comparison between the

convergence rates of different optimization algorithms. In some settings, evaluating the objective function  $f$  is itself expensive. For example, this is the case in many machine learning applications where the objective involves a sum over a large number of data points. In this case, a popular approach is to use (projected) Stochastic Gradient Descent (SGD), where it suffices to evaluate a stochastic version of the objective  $f$  at each iteration, which can often be done efficiently. SGD enjoys similar computational benefits as GD. However, when the domain  $\mathcal{C}$  of interest is not a Euclidean ball, the main bottleneck of such methods is often the projections that need be performed each time the iterates of the algorithm step outside of  $\mathcal{C}$ . This has fueled a long line of research around the design of, so-called, projection-free algorithms that swap expensive projections for cheaper linear optimizations over the domain  $\mathcal{C}$ . The reader is referred to the vast literature on projection-free optimization that include, for instance, (Frank et al., 1956; Hazan, 2008; Jaggi, 2013; Garber, 2016; Kerdreux, 2020; Bomze et al., 2021).

The first projection-free algorithm (a.k.a. the Frank-Wolfe algorithm) was introduced by Frank et al. (1956) who applied it to smooth optimization on polyhedral sets. The Frank-Wolfe algorithm swaps expensive projections on the convex set  $\mathcal{C}$  for linear optimization on  $\mathcal{C}$ , which for many sets of interest can be performed efficiently (see Table 2). A similar algorithm was later introduced for the Online Convex Optimization (OCO) setting with linear losses (Kalai and Vempala, 2005). This is a setting where instead of optimizing a fixed function, the goal is to minimize a quantity called *regret*. Formally, at each round  $t$  in OCO, a learner (algorithm) outputs an iterate  $\mathbf{x}_t$  in some convex set  $\mathcal{C}$ , then observes a convex loss function  $\ell_t : \mathcal{C} \rightarrow \mathbb{R}$  that may be chosen by an adversary based on  $\mathbf{x}_t$  and the history up to round  $t$ . The goal is to minimize the regret  $\text{Regret}_T(\mathbf{x}) := \sum_{t=1}^T \ell_t(\mathbf{x}_t) - \sum_{t=1}^T \ell_t(\mathbf{x})$ , which represents the difference between the cumulative loss of the learner and that of any comparator  $\mathbf{x} \in \mathcal{C}$ .

Offline [resp. Stochastic] optimization are special cases of OCO, where  $\ell_t$  is fixed [resp. stochastic with a fixed mean] for all  $t$ . Algorithms designed for OCO (e.g. Online Gradient Descent (Zinkevich, 2003)) often lead to optimal convergence rates when used in the offline and stochastic settings via online-to-batch conversion techniques (Cesa-Bianchi et al., 2004; Shalev-Shwartz et al., 2011; Cutkosky, 2019). For examples of problems that can be modeled via OCO we refer the reader to the OCO introduction by Hazan (2016). Unlike in the offline setting where there are algorithms such as those based on the cutting plane method that converge exponentially fast with the number of iterations, there are no equivalent options for OCO since the losses are different at each round and the goal is to minimize the regret. This means that designing efficient projection-free algorithms for OCO is even more crucial than in the offline and stochastic settings.

In this paper, we continue the long line of research in the design of efficient projection-free algorithms for online and stochastic optimization (see e.g. (Kalai and Vempala, 2005; Hazan and Kale, 2012; Hazan and Luo, 2016; Hazan and Minasyan, 2020)). We depart from the Frank-Wolfe-style approach by presenting a new algorithm that requires a Membership Oracle  $\mathcal{M}_{\mathcal{C}}$  for  $\mathcal{C}$  instead of a Linear Optimization Oracle (LOO). When given  $\mathbf{x} \in \mathbb{R}^d$ , the Membership Oracle returns  $\mathbb{I}\{\mathbf{x} \in \mathcal{C}\}$ . We present an algorithm that requires a total of  $O(T \ln T)$  calls to a (possibly approximate) Membership Oracle and guarantees a  $O(\sqrt{T})$  regret bound for general convex functions and convex sets. This regret bound is optimal in the number of rounds  $T$  for general OCO. This result improves upon previous results designed for more specific use cases, and closes the longstanding regret gap between projection-free (e.g. Frank-Wolfe-style algorithm) and projection-based methods by achieving optimal regret. We now discuss our contributions in more detail.

**Contributions.** For OCO with general convex losses, we provide an algorithm in the form of a (projection-free) reduction that achieves a  $O(\kappa\sqrt{T})$  regret after  $T$  rounds using a total of  $\tilde{O}(T)$  [resp.  $\tilde{O}(dT)$ ] calls to an approximate Separation [resp. Membership] Oracle, for some  $\kappa > 0$  that depends on  $\mathcal{C}$ . The constant  $\kappa$ , which may be assumed less than  $d$  without loss of generality (see Rem. 1), takes the role of a  $\sqrt{d}$  factor that is present in the regret bounds of previous projection-free algorithms. In Table 2, we display bounds on  $\kappa$  as function of  $d$  for many sets of interest.

We also present a more efficient version of our first algorithm that makes only  $\tilde{O}(1)$  calls per round to a Membership Oracle (instead of  $\tilde{O}(d)$  calls) at the cost of a multiplicative  $\sqrt{d}$  factor in the regret bound. Thus, our approach allows a trade-off between computation (i.e. Oracle calls) and regret. Our reduction also allows us to build an adaptive projection-free algorithm for OCO with strongly convex losses that achieves a logarithm regret without requiring the parameter of strong convexity as input, while ensuring a  $\tilde{O}(\sqrt{T})$  regret in the worst case.

Our guarantees for general OCO readily transfer to the stochastic setting via well-known online-to-batch conversion techniques (Cesa-Bianchi et al., 2004; Shalev-Shwartz et al., 2011), leading to a convergence rate of order  $O(\kappa/\sqrt{T})$  which is optimal in  $T$ . We also present an algorithm for  $\beta$ -smooth stochastic optimization that achieves the rate  $O(\beta\kappa^2/T^2 + \sigma\kappa/\sqrt{T})$ , where  $\sigma^2$  is the variance of the noise of the observed subgradients. Crucially, our algorithm does not require knowledge of either  $\beta$  or  $\sigma$ , and is thus fully adaptive. When instantiated in the offline setting (i.e. when  $\sigma = 0$ ), these rates imply that our algorithms represent viable alternatives to state-of-the-art cutting plane methods (Lee et al., 2018) whenever  $d \geq \Omega(\kappa^2/\epsilon^2)$  for general convex functions or when  $d \geq \Omega(\kappa/\sqrt{\epsilon})$  for smooth functions. We summarize our contributions in Table 1, where we compare our average regret/rates to the best known methods that apply to general convex sets.

We achieve these results by reducing an OCO problem on a general convex set  $\mathcal{C}$  to an OCO problem on a ball, where projections can be performed efficiently. More specifically, we present an algorithm wrapper that takes in any base algorithm A defined on a Euclidean ball (where projections can be performed cheaply) containing  $\mathcal{C}$  and turns it into an algorithm on  $\mathcal{C}$  that does not incur any expensive projections and whose regret bound is at most a factor  $O(\kappa)$  worse than that of A.

**Related Works.** This paper continues a long a line of work in projection-free online learning. Table 1 summarizes the guarantees of state-of-the-art Frank-Wolfe-style algorithms that are applicable to general convex sets. For OCO with general convex sets/functions, the best known regret is of order  $O(T^{3/4})$  and was achieved by Hazan and Kale (2012). The corresponding Frank-Wolfe algorithm requires a single call to a Linear Optimization Oracle per round. Later, Garber and Hazan (2016) introduced an algorithm that achieves the optimal  $O(\sqrt{T})$  [resp.  $O(\log T)$ ] regret in OCO with general [resp. strongly convex] losses with one call to an LOO per round. However, their results only apply when the convex set is a polytope, and they use an Oracle that requires at least  $O(\min(dT, d^2))$  arithmetic operations per round (translating into  $O(\min(n^4, Tn^2))$  operations when  $\mathcal{C}$  is a subset of  $n \times n$  matrices). Closer to our approach is that of Mahdavi et al. (2012) and Levy and Krause (2019) who essentially also considered Membership and Separation Oracles for their algorithms instead of a Linear Optimization Oracle on  $\mathcal{C}$ . In their setting, the set  $\mathcal{C}$  is of the form  $\{\mathbf{x} \in \mathbb{R}^d \mid h(\mathbf{x}) \leq 0\}$ , where  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  is a smooth convex function, which translates into the set  $\mathcal{C}$  being smooth. The approach of Mahdavi et al. (2012) requires at least one orthogonal projection onto  $\mathcal{C}$ . Levy and Krause (2019) avoid this by providing a fast approximate projection and achieve optimal regret bounds for both convex and strongly convex losses while only requiring a single call to a Separating and Membership Oracle for  $\mathcal{C}$  per round. However, their guarantees only applies for smooth sets.

Setting	Loss Function	Average Regret/Rate Previous Best	This paper
Online	Non-Smooth	$T^{-1/4}$ (Hazan and Kale, 2012)	$T^{-1/2}$ (Thms. 8,15, 16)
Online	Smooth	$T^{-1/3}$ (Hazan and Minasyan, 2020)	$T^{-1/2}$ (Thm. 15)
Online	Strongly Convex	$T^{-1/3}$ (Kretzu and Garber, 2021)	$T^{-1} \cdot \ln T$ (Thm. 17)
Stochastic	Non-Smooth	$T^{-1/3}$ (Hazan and Kale, 2012)	$T^{-1/2}$ (Thms. 8,15,16)
Stochastic	$\beta$ -Smooth ( $\sigma^2$ -var)	$T^{-1/2}$ (Lan et al., 2017)	$\beta T^{-2} + \sigma/T^{-1/2}$ (Thm. 18)
Stochastic	Strongly Convex	$T^{-1/3}$ (Hazan and Kale, 2012)	$T^{-1} \cdot \ln T$ (Thm. 17)
Offline	Non-Smooth	$T^{-1/3}$ (Hazan and Kale, 2012)	$T^{-1/2}$ (Thms. 8,15,16)
Offline	$\beta$ -Smooth	$\beta T^{-2}$ (Lan and Zhou, 2016)	$\beta T^{-2}$ (Thm. 18)

Table 1: Only results that hold for general bounded convex sets are included. The algorithms that achieve the reported results for our paper require either  $\tilde{O}(d)$  or  $\tilde{O}(1)$  calls per round to a Membership Oracle for  $\mathcal{C}$ . In the latter case, the regret bound is a  $\sqrt{d}$  factor worse.

In a way, their algorithm complements that of Garber and Hazan (2016) that works for polytopes (which are non-smooth sets). Our approach can be viewed as a generalization of that of Mahdavi et al. (2012); Levy and Krause (2019) to general sets. Finally, we mention the work of Abernethy et al. (2012) who used a self-concordant barrier for the set of interest to avoid projections. The efficient version of their algorithm, which is based on the damped Newton step, has a  $O(d^3)$  per-round computational complexity (due to matrix inversion).

For OCO with smooth [resp. strongly] convex losses (Hazan and Minasyan, 2020) [resp. (Kretzu and Garber, 2021)] provide algorithms that achieve a  $O(T^{2/3})$  regret (up to multiplicative factors in the dimension), requiring a single call to an LOO per round. For smooth projection-free stochastic optimization, the best rate we are aware of is of order  $O(\beta/T^2 + \sigma/\sqrt{T})$ , where  $\beta$  is the smoothness parameter and  $\sigma$  is the stochastic noise. The first algorithm to achieve this makes  $O(T)$  LOO calls after  $T$  rounds and is due to Lan et al. (2017). This improves on the previous best  $T^{-1/4}$  rate due to Hazan and Kale (2012) that uses the same number of Oracle calls. We achieve a similar rate as that of Lan et al. (2017), albeit we use a Membership Oracle instead of an LOO. The story is similar for offline smooth convex optimization where we achieve the optimal accelerated rate of  $O(\beta/T^2)$  for  $\beta$ -smooth function as the algorithm of Lan and Zhou (2016), albeit without requiring knowledge of  $\beta$ .

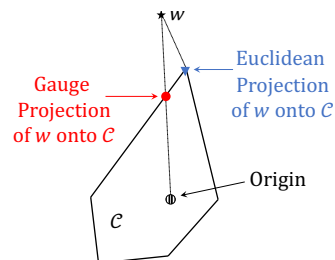


Figure 1: Gauge vs Euclidean projections.

Finally, the core idea behind our algorithm is inspired by recent techniques (Cutkosky and Orabona, 2018; Cutkosky, 2020) that leverage certain surrogate losses to reduce a constrained OCO problem to an unconstrained one. The key novelty in our work is a choice of surrogate losses that depend on the Gauge function (also known as the Minkowski function) of the set  $\mathcal{C}$ . In particular, we choose this dependence in a way that allows us to replace potentially expensive (Euclidean) projections by inexpensive Gauge projections (see Definition 5 and Figure 1), provided one has access to a Membership Oracle. We are able to perform efficient Gauge projections using recent techniques for estimating subgradients of convex functions that are not necessarily differentiable

(Lee et al., 2018). Our algorithm for strongly-convex OCO [resp. smooth stochastic optimization] builds on reductions due to Cutkosky and Orabona (2018) [resp. Cutkosky (2020)].

**Outline.** In §2, we describe our setting and provide some standard definitions from convex analysis. In §3, we describe tools we use to approximate the Gauge function of a set and its subgradients; these will be key to deriving our projection-free algorithms. In §4, we present our main method that reduces constrained OCO on an arbitrary compact convex set  $\mathcal{C}$  to OCO on a Euclidean ball containing  $\mathcal{C}$ . There, we also discuss the implications of this result in the stochastic and offline optimization settings. We conclude with a discussion in §5. A table of contents is provided in the appendix.

**Notation.** For a real function  $f: \mathcal{A} \rightarrow \mathbb{R}$ , we denote by  $\arg \max_{a \in \mathcal{A}} f(a)$  [resp.  $\arg \min_{a \in \mathcal{A}} f(a)$ ] the subset  $\mathcal{K}$  of  $\mathcal{A}$  such that  $f(b) = \sup_{a \in \mathcal{A}} f(a)$  [resp.  $f(b) = \inf_{a \in \mathcal{A}} f(a)$ ], for all  $b \in \mathcal{K}$ . We denote by  $\mathcal{B}(r)$  the Euclidean ball of radius  $r$  centered at the origin, and let  $\|\cdot\|$  be the Euclidean norm. Finally, for any mathematical statement  $\mathcal{E}$ , we let  $\mathbb{I}_{\mathcal{E}} = 1$  if  $\mathcal{E}$  is true, and 0 otherwise. Additional notation needed in our proofs is included in Appendix A.

## 2. Setting and Definitions

We consider the standard OCO setting where at each *round*  $t \geq 1$ , a learner (the algorithm) outputs an iterate  $\mathbf{x}_t$  in some convex set  $\mathcal{C} \subset \mathbb{R}^d$ , then the environment reveals a convex loss  $\ell_t: \mathcal{C} \rightarrow \mathbb{R}$ , and the learner suffers loss  $\ell_t(\mathbf{x}_t)$ . The goal of the learner is to minimize the *regret*  $\text{Regret}_T(\mathbf{x}) := \sum_{t=1}^T \ell_t(\mathbf{x}_t) - \sum_{t=1}^T \ell_t(\mathbf{x})$  after  $T \geq 1$  rounds against any fixed comparator  $\mathbf{x} \in \mathcal{C}$ .

We say that  $(\ell_t)$  is an *adversarial loss sequence* when  $\ell_t, t \geq 1$ , may depend on the learner’s iterate  $\mathbf{x}_t$  as well as the history  $(\mathbf{x}_s, \ell_s)_{s \in [t-1]}$ . This is the type of losses we will consider in our general OCO setting. We note that our bounds are often on the *linearized regret*  $\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle$ , where, for all  $t \geq 1$ ,  $\mathbf{g}_t$  is in the subdifferential  $\partial \ell_t(\mathbf{x}_t)$  of  $\ell_t$  at  $\mathbf{x}_t$  (see e.g. (Hiriart-Urruty and Lemaréchal, 2004) for a definition), and  $\langle \cdot, \cdot \rangle$  denotes the standard inner product. Such bounds automatically transfer to bounds on the regret  $\text{Regret}_T(\mathbf{x})$  by convexity of the losses  $(\ell_t)$ ; for any  $\mathbf{x} \in \mathcal{C}$  and  $t \geq 1$ , we have  $\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{x}) \leq \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle$ , for all  $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$ . Bounding the linearized regret is standard in OCO (see e.g. (Hazan, 2019)).

In this work, we will not assume knowledge of the horizon  $T$ , and so our regret bounds hold for all  $T$ ’s simultaneously; these are so-called *anytime* regret bound. We will allow the outputs  $(\mathbf{x}_t)$  of the learner to be random. In this case, for any  $t \geq 1$ , we let  $\mathcal{G}_t$  denote the  $\sigma$ -algebra generated by  $(\mathbf{x}_s)_{s \in [t]}$ , and we denote  $\mathbb{E}_t[\cdot] := \mathbb{E}[\cdot \mid \mathcal{G}_t]$ . When, for any round  $t$ , the output  $\mathbf{x}_t$  of an algorithm A is a deterministic function of the history  $(\mathbf{x}_s, \ell_s)_{s \in [t-1]}$ , we say that A is *deterministic*. Throughout, we assume that  $\mathcal{C}$  is “sandwiched” between two Euclidean balls:

**Assumption 1** We assume that  $\mathcal{C} \subseteq \mathbb{R}^d$  is a closed convex set s.t. for some  $0 < r \leq R$ , we have

$$\mathcal{B}(r) \subseteq \mathcal{C} \subseteq \mathcal{B}(R). \quad (1)$$

**Remark 1** Assumption 1 comes with no loss of generality as one can easily re-parametrize any OCO setting to satisfy (1) without any significant computational overhead. In Appendix F, we discuss the details of how this can be done in the general OCO setting as well as in various popular optimization settings. In Appendix F, we also derive corresponding upper bounds on the *asphericity*  $\kappa := R/r$  (Goffin, 1988), which appears as a multiplicative factor in our regret bounds. As can be



seen in Table 2,  $\kappa$  is less than  $\sqrt{d}$  for many sets of interest. We also note that, in the worst case,  $\kappa$  can always be bounded by  $d$  [resp.  $\sqrt{d}$ ] for general [resp. centrally symmetric] convex sets by applying an affine re-parametrization (Flaxman et al., 2005).

Our goal in this paper is to design algorithms with low regret in OCO without incurring expensive (Euclidean) projections onto the set  $\mathcal{C}$ . For this, we will use a Membership Oracle for the set  $\mathcal{C}$  to perform what we call Gauge Projections, which can be done efficiently using a small number of calls to the Oracle. A Membership Oracle is a map  $\mathcal{M}_{\mathcal{C}}: \mathbb{R}^d \rightarrow \{0, 1\}$ , where for any  $\mathbf{x} \in \mathbb{R}^d$ ,  $\mathcal{M}_{\mathcal{C}}(\mathbf{x}) = 1$  if and only if  $\mathbf{x} \in \mathcal{C}$ . Such Oracles are often the building blocks of many offline optimization algorithms (Grötschel et al., 1993; Bubeck, 2015; Lee et al., 2018). In many practical settings, one can only efficiently test membership approximately, and so for the sake of generality we will state our results when only such approximate Oracles are available. To formally define the notion of approximate Membership Oracle, we let  $\mathcal{B}(\mathcal{C}, \delta) := \{\mathbf{x} \in \mathbb{R}^d : \exists \mathbf{y} \in \mathcal{C} \text{ such that } \|\mathbf{x} - \mathbf{y}\| \leq \delta\}$ . We also define  $\mathcal{B}(\mathcal{C}, -\delta) := \{\mathbf{x} \in \mathbb{R}^d : \forall \mathbf{y} \in \mathcal{B}(\delta), \mathbf{x} + \mathbf{y} \in \mathcal{C}\}$ .

**Definition 2 (Approximate Membership Oracle)** *Let  $\delta > 0$  and  $\mathcal{C}$  be a convex set s.t.  $\mathcal{B}(\delta) \subseteq \mathcal{C} \subseteq \mathbb{R}^d$ . Then, the map  $\text{MEM}_{\mathcal{C}}(\cdot; \delta): \mathbb{R}^d \rightarrow \{0, 1\}$  is a  $\delta$ -approximate Membership Oracle if*

$$\forall \mathbf{x} \in \mathbb{R}^d, \quad \mathbb{I}_{\{\mathbf{x} \in \mathcal{B}(\mathcal{C}, -\delta)\}} \wedge \mathbb{I}_{\{\mathbf{x} \in \mathcal{B}(\mathcal{C}, \delta)\}} \leq \text{MEM}_{\mathcal{C}}(\mathbf{x}; \delta) \leq \mathbb{I}_{\{\mathbf{x} \in \mathcal{B}(\mathcal{C}, -\delta)\}} \vee \mathbb{I}_{\{\mathbf{x} \in \mathcal{B}(\mathcal{C}, \delta)\}}. \quad (2)$$

Some definitions of (approximate) Membership Oracles only require (2) to hold with high probability (see e.g. Lee et al. (2018)). Our analysis can easily be extended to this case, but we make the deterministic assumption in (2) to simplify the presentation. We will use an approximate Membership Oracle to build a  $\delta$ -approximate Linear Optimization Oracle  $\mathcal{O}_{\mathcal{C}^\circ}(\cdot; \delta)$  for the polar set  $\mathcal{C}^\circ$ , where for any  $\mathbf{w}$  in a ball containing  $\mathcal{C}$ , the random output  $(\gamma, \mathbf{s})$  of  $\mathcal{O}_{\mathcal{C}^\circ}(\mathbf{w}; \delta)$  essentially satisfies  $\sup_{\mathbf{u} \in \mathcal{C}} \langle \mathbf{w}, \mathbf{u} \rangle \leq \gamma \leq \sup_{\mathbf{u} \in \mathcal{C}} \langle \mathbf{w}, \mathbf{u} \rangle + \delta$  and  $\langle \mathbf{s}, \mathbf{w} \rangle \geq \sup_{\mathbf{u} \in \mathcal{C}} \langle \mathbf{w}, \mathbf{u} \rangle - \delta$  with high probability<sup>1</sup>.

We now present some standard definitions from convex analysis:

**Definition 3 (Polar Set)** *Let  $\mathcal{C} \subseteq \mathbb{R}^d$  be a convex set. The polar set  $\mathcal{C}^\circ$  of  $\mathcal{C}$  is defined as  $\mathcal{C}^\circ := \{\mathbf{s} \in \mathbb{R}^d : \langle \mathbf{s}, \mathbf{u} \rangle \leq 1, \forall \mathbf{u} \in \mathcal{C}\}$ .*

The notation of a polar set is key to our approach. We will essentially build an efficient Linear Optimization Oracle on  $\mathcal{C}^\circ$ , which will then allow us to perform, what we call, Gauge projections instead of Euclidean ones. Gauge projections (Def. 5 below) are defined with respect to a quasi-metric based on the *Gauge function* (also known as the Minkowski function) of the set  $\mathcal{C}$ :

**Definition 4 (Gauge and Support Functions)** *Let  $\mathcal{C} \subseteq \mathbb{R}^d$  be a convex set. The Gauge function  $\gamma_{\mathcal{C}}$  of  $\mathcal{C}$  is defined as  $\gamma_{\mathcal{C}}(\mathbf{w}) := \inf\{\lambda > 0 : \mathbf{w} \in \lambda\mathcal{C}\}$ , for  $\mathbf{w} \in \mathbb{R}^d$ . The support function  $\sigma_{\mathcal{C}}$  of  $\mathcal{C}$  is defined as  $\sigma_{\mathcal{C}}(\mathbf{w}) := \sup_{\mathbf{u} \in \mathcal{C}} \langle \mathbf{u}, \mathbf{w} \rangle$ , for  $\mathbf{w} \in \mathbb{R}^d$ .*

The Gauge function satisfies all properties of a norm except for symmetry, which requires the set  $\mathcal{C}$  to be centrally symmetric. Fortunately, we do not need this property to perform efficient projections with respect to the quasi-metric  $d_{\mathcal{C}}(\mathbf{u}, \mathbf{w}) := \gamma_{\mathcal{C}}(\mathbf{u} - \mathbf{w})$ . The concept of support function in the previous definition will also be key in our analysis (it is the Gauge function of the polar set  $\mathcal{C}^\circ$ ). We now introduce the concept of Gauge projection/distance to a set:

1. The guarantee provided by  $\gamma$  (i.e.  $\langle \mathbf{w}, \mathbf{u} \rangle \leq \gamma \leq \sup_{\mathbf{u} \in \mathcal{C}} \langle \mathbf{w}, \mathbf{u} \rangle + \delta$ ) will be crucial in our analysis, which is why we make our LO Oracle  $\mathcal{O}_{\mathcal{C}^\circ}$  return a scalar-vector pair instead of a single vector as is more common for such Oracles.

**Definition 5** Let  $\mathcal{C}$  be as in Assump. 1, the Gauge projection on  $\mathcal{C}$  is the set-valued map  $\Pi_{\mathcal{C}}^{\text{G}} : \mathbb{R}^d \rightrightarrows \mathcal{C}$  defined by  $\Pi_{\mathcal{C}}^{\text{G}}(\mathbf{w}) := \arg \min_{\mathbf{u} \in \mathcal{C}} \gamma_{\mathcal{C}}(\mathbf{w} - \mathbf{u})$ . We also define the Gauge distance to  $\mathcal{C}$  as the map:

$$S_{\mathcal{C}}(\mathbf{w}) = \inf_{\mathbf{u} \in \mathcal{C}} \gamma_{\mathcal{C}}(\mathbf{w} - \mathbf{u}) \geq 0, \quad \text{for all } \mathbf{w} \in \mathbb{R}^d. \quad (3)$$

Gauge projections, which can be done efficiently using a Membership Oracle (see §3), will be all we need to ensure the iterates of our algorithms are in  $\mathcal{C}$ .

### 3. Preliminaries: Gauge Function Approximation

In this section, we briefly describe some of the technical tools we require for our method (additional details are provided in Appendix B). At the heart of our approach is the design of surrogate losses for online and stochastic optimization that depend on the Gauge distance function  $S_{\mathcal{C}}$  in (3) in a way that leads to algorithms that perform Gauge projections instead of (potentially expensive) Euclidean projections. Naturally, this will require the ability to efficiently evaluate  $S_{\mathcal{C}}$  and its subgradients. The next lemma (whose proof is in App. B) shows that  $S_{\mathcal{C}}$  and  $\partial S_{\mathcal{C}}$  can be expressed as concise functions of  $\gamma_{\mathcal{C}}$  (the Gauge function):

**Lemma 6** For any set  $\mathcal{C}$  satisfying Assump. 1, and any  $\mathbf{u} \notin \mathcal{C}$  and  $\mathbf{w} \in \mathbb{R}^d$ , we have  $\frac{\mathbf{u}}{\gamma_{\mathcal{C}}(\mathbf{u})} \in \Pi_{\mathcal{C}}^{\text{G}}(\mathbf{u})$ ,

$$S_{\mathcal{C}}(\mathbf{w}) = 0 \vee (\gamma_{\mathcal{C}}(\mathbf{w}) - 1), \quad \text{and} \quad \partial S_{\mathcal{C}}(\mathbf{w}) = \begin{cases} \partial \gamma_{\mathcal{C}}(\mathbf{w}) = \arg \max_{\mathbf{s} \in \mathcal{C}^{\circ}} \langle \mathbf{s}, \mathbf{w} \rangle, & \text{if } \mathbf{w} \notin \mathcal{C}; \\ \{\mathbf{0}\}, & \text{otherwise.} \end{cases}$$

Note that for any  $\mathbf{u} \notin \mathcal{C}$ , Lemma 6 also provides the expression of a Gauge projection point in  $\Pi_{\mathcal{C}}^{\text{G}}(\mathbf{u})$  (i.e.  $\mathbf{u}/\gamma_{\mathcal{C}}(\mathbf{u})$ ), which will be leveraged by our algorithm. Next, we describe how we build efficient algorithms for approximating  $\gamma_{\mathcal{C}}$  and its subgradients using our approximate Membership Oracle  $\text{MEM}_{\mathcal{C}}$ . As Lemma 6 shows, this is all we need to approximate  $S_{\mathcal{C}}$ , its subgradients, and Gauge projections onto  $\mathcal{C}$ .

**Approximating  $\gamma_{\mathcal{C}}$ .** By definition of the Gauge function, we have,

$$\forall \mathbf{w} \in \mathbb{R}^d, \quad \gamma_{\mathcal{C}}(\mathbf{w}) = \inf\{\lambda \in \mathbb{R}_{\geq 0} \mid \mathbf{w} \in \lambda \mathcal{C}\} = 1 / \sup\{\nu \in \mathbb{R}_{\geq 0} \mid \nu \mathbf{w} \in \mathcal{C}\}. \quad (4)$$

Using the Membership Oracle  $\text{MEM}_{\mathcal{C}}$ , we can approximate the largest  $\nu \geq 0$  such that  $\nu \mathbf{w} \in \mathcal{C}$  via bisection, which will lead to an approximation of  $\gamma_{\mathcal{C}}(\mathbf{w})$  by (4). This is exactly what we do in Algorithm 2 in Appendix B. Algorithm 2 ( $\text{GAU}_{\mathcal{C}}$ ) requires at most  $\lceil \log_2((4\kappa)^2/\delta) \rceil + 1$  calls to the Membership Oracle  $\text{MEM}_{\mathcal{C}}(\cdot; r\delta/(4\kappa)^2)$  to return a  $\delta$ -approximate value of  $\gamma_{\mathcal{C}}$ , where  $\kappa := R/r$  and  $r, R > 0$  are as in (1). In Lemma 10 of Appendix B, we state the full guarantee of Algorithm 2.

**Approximating the subgradients of  $\gamma_{\mathcal{C}}$ .** In addition to approximating  $\gamma_{\mathcal{C}}$ , we will also need to approximate its subgradients. By Lemma 6, the subdifferential of  $S_{\mathcal{C}}$  coincides with that of  $\gamma_{\mathcal{C}}$  at any point  $\mathbf{w}$  outside  $\mathcal{C}$ . The lemma also shows that finding a subgradient of  $\gamma_{\mathcal{C}}$  is equivalent to performing linear optimization on the polar set  $\mathcal{C}^{\circ}$ . In Appendix B, we present two algorithms,  $\text{OPT}_{\mathcal{C}^{\circ}}$  (Alg. 3) and  $\text{OPT}_{\text{Id}, \mathcal{C}^{\circ}}$  (Alg. 4), based on (Lee et al., 2018, Algorithm 2), which use  $\text{GAU}_{\mathcal{C}}$  (Alg. 2) and a random partial difference along different coordinates to approximate the subgradients of  $\gamma_{\mathcal{C}}$ . The first algorithm  $\text{OPT}_{\mathcal{C}^{\circ}}$  requires at most  $O(d \ln(d\kappa/\delta))$  calls to the approximate Membership Oracle  $\text{MEM}_{\mathcal{C}}$  to find a “ $\delta$ -approximate” subgradient; the precise guarantee is stated in Proposition 11.

The algorithm  $\text{OPT}_{1d, \mathcal{C}^\circ}$  is a stochastic version of  $\text{OPT}_{\mathcal{C}^\circ}$  that picks a single random coordinate  $I$  along which to estimate the subgradient of the Gauge function  $\gamma_{\mathcal{C}}$ . As a result,  $\text{OPT}_{1d, \mathcal{C}^\circ}$  requires at most  $O(\ln(d\kappa/\delta))$  calls to the approximate Membership Oracle  $\text{MEM}_{\mathcal{C}}$ , which will provide a computational benefit over  $\text{OPT}_{\mathcal{C}^\circ}$  at the cost of a multiplicative  $\sqrt{d}$  in the regret bound. Further, the output of  $\text{OPT}_{1d, \mathcal{C}^\circ}$  is equal to that of  $\text{OPT}_{\mathcal{C}^\circ}$  in expectation for the same input (see Lemma 13).

We now move to the main section of the paper where we present our projection-free reduction that uses the tools we described in this section.

#### 4. Efficient Projection-Free Online and Stochastic Convex Optimization

In this section, we will use our approximate Linear Optimization Oracles on  $\mathcal{C}^\circ$  described in the previous section to build efficient algorithms for OCO with optimal regret bounds. Our final algorithms make at most  $\tilde{O}(T)$  calls to the Membership Oracle  $\text{MEM}_{\mathcal{C}}$  after  $T \geq 1$  rounds, and are also time-uniform in the sense that they do not require a horizon  $T$  as input.

---

**Algorithm 1** Projection-Free Algorithm-Wrapper for OCO on  $\mathcal{C}$ .

---

**Require:** **I**) An OCO algorithm  $A$  on  $\mathcal{B}(R) \supseteq \mathcal{C}$ ; **II**) A tolerance sequence  $(\delta_t) \subset (0, 1/3)$ ; and **III**) LO Oracle  $\mathcal{O}_{\mathcal{C}^\circ}$  on  $\mathcal{C}^\circ$  such that  $\mathcal{O}_{\mathcal{C}^\circ}(\mathbf{w}; \delta) \in \mathbb{R}_{\geq 0} \times \mathbb{R}^d, \forall \mathbf{w} \in \mathcal{B}(R), \delta \in (0, 1)$ .  
*// In the ideal case,  $(\gamma, \mathbf{s}) = \mathcal{O}_{\mathcal{C}^\circ}(\mathbf{w}; \delta) \implies \gamma = \langle \mathbf{s}, \mathbf{w} \rangle$  and  $\mathbf{s} \in \arg \max_{\mathbf{u} \in \mathcal{C}^\circ} \langle \mathbf{u}, \mathbf{w} \rangle$*

- 1: Initialize  $A$  and set  $\mathbf{w}_1 \in \mathcal{B}(R)$  to  $A$ 's first output.
- 2: **for**  $t = 1, 2, \dots$  **do**
- 3:   Set  $(\gamma_t, \mathbf{s}_t) = \mathcal{O}_{\mathcal{C}^\circ}(\mathbf{w}_t; \delta_t)$
- 4:   Set  $\mathbf{v}_t = \mathbb{I}_{\{\gamma_t \geq 1\}} \mathbf{s}_t$ . *// Subgradient of  $S_{\mathcal{C}}$  at  $\mathbf{w}_t$*
- 5:   Play  $\mathbf{x}_t = \mathbb{I}_{\{\gamma_t \geq 1\}} \mathbf{w}_t / \gamma_t + \mathbb{I}_{\{\gamma_t < 1\}} \mathbf{w}_t$ . *// Gauge projection of  $\mathbf{w}_t$  onto  $\mathcal{C}$*
- 6:   Observe subgradient  $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$ .
- 7:   Set  $\tilde{\mathbf{g}}_t = \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \mathbf{v}_t$  *//  $\tilde{\mathbf{g}}_t \in \partial \tilde{\ell}_t(\mathbf{w}_t)$ ;  $\tilde{\ell}_t(\mathbf{w}) := \langle \mathbf{g}_t, \mathbf{w} \rangle - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle S_{\mathcal{C}}(\mathbf{w})$*
- 8:   Set  $A$ 's  $t$ th loss function to  $f_t : \mathbf{w} \mapsto \langle \tilde{\mathbf{g}}_t, \mathbf{w} \rangle$ .
- 9:   Set  $\mathbf{w}_{t+1} \in \mathcal{B}(R)$  to  $A$ 's  $(t+1)$ th output given the history  $((\mathbf{w}_i, \mathbf{x}_i, f_i)_{i \leq t})$ .
- 10: **end for**

---

To avoid expensive projections, our algorithms make use of convex surrogate losses of the form  $\tilde{\ell}(\mathbf{w}) := \langle \mathbf{g}, \mathbf{w} \rangle + b S_{\mathcal{C}}(\mathbf{w})$ , for  $\mathbf{w} \in \mathbb{R}^d$  and  $b \geq 0$ , where  $S_{\mathcal{C}}$  is the Gauge distance function in (3) and  $\mathbf{g}$  is a subgradient of the loss of interest  $\ell : \mathcal{C} \rightarrow \mathbb{R}$ . Note that  $\tilde{\ell}$  is not constrained to the set  $\mathcal{C}$ , unlike the actual loss of interest  $\ell$ . The choice of such a surrogate loss function is inspired by existing constrained-to-unconstrained reductions in OCO due to Cutkosky and Orabona (2018); Cutkosky (2020). Similar to the latter, our projection-free reduction allows us to bound the regret of our Algorithm 1 by the regret of any OCO subroutine  $A$  that is fed subgradients of the surrogate losses. Here, the iterates of  $A$  need not be constrained to the set  $\mathcal{C}$  since the domain of the surrogate losses is technically unconstrained. Thus, to build efficient projection-free OCO algorithms on  $\mathcal{C}$  with optimal regret, it suffices to pick a sub-algorithm  $A$  that I) has an optimal regret bound on a set  $\mathcal{W}$  containing  $\mathcal{C}$ , and II) does not incur any expensive projections. We will pick  $\mathcal{W} = \mathcal{B}(R) \supseteq \mathcal{C}$  so that Euclidean projections onto  $\mathcal{W}$ , which may be required by  $A$ , can be performed efficiently.

We now state our main reduction result when the LO Oracle  $\mathcal{O}_{\mathcal{C}^\circ}$  is set to the approximate LO Oracle  $\text{OPT}_{\mathcal{C}^\circ}$  (Alg. 3) we described in §3. The full proof of the lemma is in App. I.1.



**Lemma 7** Let  $\kappa := R/r$  and  $A$  be any OCO algorithm on  $\mathcal{B}(R)$ . Further, for  $t \geq 1$ , let  $\mathbf{w}_t, \tilde{\mathbf{g}}_t, \mathbf{x}_t$ , and  $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$  be as in Alg. 1 with  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{\mathcal{C}^\circ}$  (Alg. 3) and any tolerance sequence  $(\delta_s) \subset (0, 1/3)$ . Then, for all  $t \geq 1$ , we have  $\mathbf{x}_t \in \mathcal{C}$ , and there exists a random variable  $\Delta_t \in [0, 15^2 d^4 \kappa^3 \delta_t^{-2}]$  such that  $\mathbb{E}_{t-1}[\Delta_t] \leq \delta_t$ ,  $\|\tilde{\mathbf{g}}_t\| \leq (1 + \Delta_t + \kappa)\|\mathbf{g}_t\|$ , and for all  $\mathbf{x} \in \mathcal{C}$ ,

$$\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \tilde{\ell}_t(\mathbf{w}_t) - \tilde{\ell}_t(\mathbf{x}) + \delta_t R \|\mathbf{g}_t\| \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + (2\delta_t + \Delta_t)R \|\mathbf{g}_t\|, \quad (5)$$

where  $\tilde{\ell}_t(\mathbf{w}) := \langle \mathbf{g}_t, \mathbf{w} \rangle - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle S_{\mathcal{C}}(\mathbf{w})$ .

**Proof Sketch.** For simplicity, we provide a proof of the lemma when Alg. 1 has access to a perfect LO Oracle  $\mathcal{O}_{\mathcal{C}^\circ}$  on  $\mathcal{C}^\circ$ , i.e. when  $\forall \mathbf{w} \in \mathbb{R}^d, \delta > 0, (\gamma, \mathbf{s}) = \mathcal{O}_{\mathcal{C}^\circ}(\mathbf{w}; \delta)$  only if  $\mathbf{s} \in \arg \max_{\mathbf{x} \in \mathcal{C}^\circ} \langle \mathbf{x}, \mathbf{w} \rangle$  and  $\gamma = \langle \mathbf{s}, \mathbf{w} \rangle$ . Thus, we will show the claims of the lemma with  $\Delta_t = \delta_t = 0$ .

Let  $\gamma_t, \mathbf{v}_t, \mathbf{w}_t, \tilde{\mathbf{g}}_t$ , and  $\mathbf{x}_t$  be as in Algorithm 1. By the expression of the subdifferential of  $S_{\mathcal{C}}$  in Lemma 6 and the definition of  $\tilde{\ell}_t$ , we have that  $\mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \mathbf{v}_t \in \partial \tilde{\ell}_t(\mathbf{w}_t)$ . Thus, since  $-\mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle > 0$  (see definition of  $\mathbf{x}_t$  in Alg. 1),  $\tilde{\ell}_t$  is convex and so

$$\forall \mathbf{x} \in \mathcal{C}, \quad \tilde{\ell}_t(\mathbf{w}_t) - \tilde{\ell}_t(\mathbf{x}) \leq \langle \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \mathbf{v}_t, \mathbf{w}_t - \mathbf{x} \rangle = \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle.$$

It remains to show that  $\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \tilde{\ell}_t(\mathbf{w}_t) - \tilde{\ell}_t(\mathbf{x})$ , for all  $t \geq 1$  and  $\mathbf{x} \in \mathcal{C}$ . First note that for all  $\mathbf{x} \in \mathcal{C}$ , we have  $S_{\mathcal{C}}(\mathbf{x}) = 0$ , and so

$$\tilde{\ell}_t(\mathbf{x}) = \langle \mathbf{g}_t, \mathbf{x} \rangle, \quad \forall \mathbf{x} \in \mathcal{C}. \quad (6)$$

We will now compare  $\langle \mathbf{g}_t, \mathbf{x}_t \rangle$  to  $\tilde{\ell}_t(\mathbf{w}_t)$  by considering cases. Suppose that  $\gamma_t < 1$ . In this case, we have  $\mathbf{x}_t = \mathbf{w}_t$  and so  $\langle \mathbf{g}_t, \mathbf{x}_t \rangle = \langle \mathbf{g}_t, \mathbf{w}_t \rangle = \tilde{\ell}_t(\mathbf{w}_t)$ . Now suppose that  $\gamma_t \geq 1$  and  $\langle \mathbf{g}_t, \mathbf{w}_t \rangle \geq 0$ . By the fact that  $\gamma_{\mathcal{C}}(\mathbf{w}_t) = \gamma_t \geq 1$  and  $\mathbf{x}_t = \mathbf{w}_t / \gamma_t$ , we have

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle \leq \langle \mathbf{g}_t, \mathbf{w}_t \rangle = \tilde{\ell}_t(\mathbf{w}_t). \quad (7)$$

Now suppose that  $\gamma_t \geq 1$  and  $\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0$ . Using the fact that  $\mathbf{x}_t = \mathbf{w}_t / \gamma_t = \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)$  and that  $S_{\mathcal{C}}(\mathbf{w}_t) = \gamma_{\mathcal{C}}(\mathbf{w}_t) - 1$  (Lemma 6), we have

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle + \langle \mathbf{g}_t, \mathbf{x}_t \rangle \cdot S_{\mathcal{C}}(\mathbf{w}_t) = \langle \mathbf{g}_t, \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t) \rangle + \langle \mathbf{g}_t, \mathbf{w}_t \rangle - \langle \mathbf{g}_t, \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t) \rangle = \langle \mathbf{g}_t, \mathbf{w}_t \rangle.$$

Rearranging this, we get

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle = \langle \mathbf{g}_t, \mathbf{w}_t \rangle - \langle \mathbf{g}_t, \mathbf{x}_t \rangle S_{\mathcal{C}}(\mathbf{w}_t) = \tilde{\ell}_t(\mathbf{w}_t). \quad (8)$$

By combining (6), (7), and (8), we obtain  $\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \tilde{\ell}_t(\mathbf{w}_t) - \tilde{\ell}_t(\mathbf{x}) \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle$ , for all  $\mathbf{x} \in \mathcal{C}$ , which shows (5) with  $\Delta_t = \delta_t = 0$ . It remains to bound  $\|\tilde{\mathbf{g}}_t\|$  in terms of  $\|\mathbf{g}_t\|$  and show that  $\mathbf{x}_t \in \mathcal{C}$ . When  $\gamma_t < 1$ , we have  $\tilde{\mathbf{g}}_t = \mathbf{g}_t$  and  $\mathbf{x}_t = \mathbf{w}_t$ , and so  $\|\tilde{\mathbf{g}}_t\| = \|\mathbf{g}_t\|$ . Furthermore, since  $\gamma_{\mathcal{C}}(\mathbf{w}_t) = \gamma_t < 1$ , the definition of the Gauge function implies that  $\mathbf{w}_t \in \mathcal{C}$  and so  $\mathbf{x}_t \in \mathcal{C}$  (since  $\mathbf{x}_t = \mathbf{w}_t$ ). Now suppose that  $\gamma_t \geq 1$ . In this case, we have  $\mathbf{x}_t = \mathbf{w}_t / \gamma_t = \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)$  and  $\tilde{\mathbf{g}}_t = \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \mathbf{v}_t$ . Using this and that  $\gamma_{\mathcal{C}}(\mathbf{w}_t) = \gamma_t \geq 1$ , we get that  $\|\tilde{\mathbf{g}}_t\|$  is equal to

$$\|\mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \mathbf{v}_t\| = \|\mathbf{g}_t\| + \|\mathbf{g}_t\| \frac{\|\mathbf{w}_t\|}{\gamma_{\mathcal{C}}(\mathbf{w}_t)} \|\mathbf{v}_t\| \stackrel{(*)}{\leq} \|\mathbf{g}_t\| \left( 1 + \frac{\|\mathbf{w}_t\|}{r} \right) \leq (1 + \kappa) \|\mathbf{g}_t\|.$$

where (\*) follows by the fact that  $\|\mathbf{v}_t\| \leq 1/r$  (since  $\mathbf{v}_t \in \mathcal{C}^\circ \subseteq \mathcal{B}(1/r)$ —see Lem. 42 for the set inclusion). Further, when  $\gamma_t \geq 1$ , we have  $\gamma_{\mathcal{C}}(\mathbf{x}_t) = \gamma_{\mathcal{C}}(\mathbf{w}_t/\gamma_{\mathcal{C}}(\mathbf{w}_t)) = \gamma_{\mathcal{C}}(\mathbf{w}_t)/\gamma_{\mathcal{C}}(\mathbf{w}_t) = 1$ , and so  $\mathbf{x}_t \in \mathcal{C}$ .  $\blacksquare$

Lemma 7 allows us to bound the instantaneous (linearized) regret  $\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle$  of Alg. 1 w.r.t. the instantaneous regret  $\langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle$  of subroutine A. Thus by summing (5) in Lemma 7 for  $t = 1, 2, \dots$ , we can bound the regret of Alg. 1 with respect to that of A. Further, Lemma 7 shows that the scale of the subgradients ( $\tilde{\mathbf{g}}_t$ ) of the losses fed to A is at most a constant times the scale of the subgradients ( $\mathbf{g}_t$ ) at the iterates ( $\mathbf{x}_t$ ) of Alg. 1; in particular,  $\|\tilde{\mathbf{g}}_t\| \leq (1 + \Delta_t + \kappa)\|\mathbf{g}_t\|$ ,  $\forall t \geq 1$ , where  $\kappa := R/r$ . This latter fact ensures that the regret bound of A is not too large as a function of the scale of ( $\mathbf{g}_t$ ).

#### 4.1. Algorithm for General Online Convex Optimization

To streamline the discussions that follow, we will instantiate Alg. 1 and state its regret bound with a specific subroutine A (for a general A, see App. D). In particular, we will set A to Follow-The-Regularized-Leader Proximal (FTRL-prox) (McMahan, 2017). FTRL-prox is summarized in Alg. 8 in App. G. The algorithm takes input parameter  $R > 0$  and performs at most one Euclidean projection onto the ball  $\mathcal{B}(R)$  per iteration. An advantage of FTRL-prox is that it has a regret bound (see Prop. 25) of the form  $O(\sqrt{V_T})$ , where  $V_T := \sum_{t=1}^T \|\mathbf{g}_t\|^2$  and ( $\mathbf{g}_t$ ) are the observed subgradients. This is desirable for reasons we outline in §4.2. We now state the regret bound of Alg. 1 when A is set to FTRL-prox:

**Theorem 8** *Let  $\delta \in (0, 1/3)$  and  $\kappa := R/r$ , with  $r$  and  $R$  as in (1). Let  $(\ell_t)$  be any adversarial sequence of convex losses on  $\mathcal{C}$  and  $(\mathbf{x}_t)$  be the iterates of Alg. 1 in response to  $(\ell_t)$ . If Alg. 1 is run with subroutine A set to FTRL-prox with parameter  $R$ ;  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{\mathcal{C}^\circ}$  (Alg. 3); and  $\delta_t = \delta/t^2$ ,  $t \geq 1$ , then for all  $T \geq 1$  and  $\rho \in (0, 1)$ , we have  $(\mathbf{x}_t)_{t \in [T]} \subset \mathcal{C}$  and with probability at least  $1 - \rho$ :*

$$\forall \mathbf{x} \in \mathcal{C}, \quad \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq 4(1 + \kappa)R\sqrt{V_T} + R \cdot (12 + 10\delta/\rho) \cdot \max_{t \leq T} \|\mathbf{g}_t\|, \quad (9)$$

where  $V_T := \sum_{t=1}^T \|\mathbf{g}_t\|^2$  and  $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$ , for all  $t \geq 1$ .

**Proof Sketch.** By Lemma 7, we have  $\mathbf{x}_t \in \mathcal{C}$ ,  $\forall t \geq 1$ . Lemma 7 also says that there exists a sequence  $(\Delta_t) \subset \mathbb{R}_{\geq 0}$  of non-negative random variables satisfying  $\mathbb{E}_{t-1}[\Delta_t] \leq \delta_t$ , for all  $t \in [T]$ , such that for all  $\mathbf{x}$  and  $t \geq 1$ ,  $\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + (2\delta_t + \Delta_t)R\|\mathbf{g}_t\|$ . Letting  $Q_t := \sum_{s=1}^t \|\tilde{\mathbf{g}}_s\|^2$ , and summing the latter inequality for  $t = 1, \dots, T$ , we obtain,

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle - \sum_{t=1}^T (2\delta_t + \Delta_t)R\|\mathbf{g}_t\| \leq \sum_{t=1}^T \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle \leq 2R\sqrt{2Q_T}, \quad \text{for all } \mathbf{x} \in \mathcal{C}, \quad (10)$$

where the last inequality follows by the fact that A is FTRL-prox and the regret bound of the latter in Prop. 25. Now, by the fact that  $\|\tilde{\mathbf{g}}_t\| \leq (1 + \kappa + \Delta_t)\|\mathbf{g}_t\|$  (Lem. 7) and  $(a + b)^2 \leq 2a^2 + 2b^2$  for all  $a, b \in \mathbb{R}_{>0}$ , the RHS of (10) can be further bounded by  $4(1 + \kappa)R\sqrt{V_T} + 4\sum_{t=1}^T \Delta_t R\|\mathbf{g}_t\|$ . It remains to bound the sum  $\sum_{t=1}^T (2\delta_t + 5\Delta_t)R\|\mathbf{g}_t\|$  in terms of  $\delta R \max_{t \leq T} \|\mathbf{g}_t\|/\rho$ , which we do using Doob's martingale inequality (see App. I.2 for the full proof).  $\blacksquare$

The  $O(\sqrt{V_T})$  regret bound in Theorem 8 is known as an *adaptive* regret bound, and can be much smaller than the standard  $O(\sqrt{T})$  worst-case regret; e.g. when the losses are smooth (Srebro et al., 2010a; Cutkosky and Orabona, 2018). We say more on this in §4.2 below.

**Oracle Complexity.** Note that at each round  $t$ , the instance of Alg. 1 in Thm. 8 invokes  $\text{OPT}_{\mathcal{C}^\circ}(\cdot; \delta_t)$  once. Thus, in light of the discussion in §3 (see also Rem. 12 in App. B), the algorithm makes at most  $O(dT \ln(dT\kappa/\delta))$  calls to the Membership Oracle  $\text{MEM}_{\mathcal{C}}$  after  $T$  rounds. We also note that if one has a Separation Oracle for  $\mathcal{C}$ , then Alg. 1 can be executed with only  $\tilde{O}(1)$  calls-per-round to the latter (instead of  $\tilde{O}(d)$  calls to  $\text{MEM}_{\mathcal{C}}$ ) without compromising the regret bound in Thm. 8.

We can reduce the number of  $\text{MEM}_{\mathcal{C}}$  Oracle calls to a total of  $O(T \ln(dT\kappa/\delta))$  by using the “one-dimensional” stochastic version of  $\text{OPT}_{\mathcal{C}^\circ}$  we discussed in §3; i.e.  $\text{OPT}_{1d, \mathcal{C}^\circ}$  (Alg. 4). This will come at the cost of a  $O(\sqrt{d})$  multiplicative factor in the regret bound (see App. C.1).

**The strongly convex case.** In App. C.2, we design a subroutine A (Alg. 5) such that the instance of Alg. 1 with this subroutine achieves a logarithmic regret whenever the losses are strongly convex, while maintaining a  $\tilde{O}(\sqrt{T})$  regret in the worse case. Further, the algorithm does not require any curvature parameter as input. Finally, in Appendix E we present a general reduction that makes our algorithm for the strongly convex case as well as many previous online algorithms (such as those by (Ross et al., 2013; Wintenberger, 2017; Kotłowski, 2017; Kempka et al., 2019)) scale-invariant; in the sense that scaling the losses by some positive constant does not change the outputs of the algorithm (consequently the algorithm does not require any scale information as input). The reduction is a generalization of previous algorithms by Mhammedi et al. (2019); Mhammedi and Koolen (2020).

## 4.2. Algorithms for Stochastic and Offline Optimization

**The stochastic setting.** The guarantee of Alg. 1 in Thm. 8 readily transfers to the stochastic setting (formally described in App. C.3) via well-known online-to-batch conversion techniques (Cesa-Bianchi et al., 2004; Shalev-Shwartz et al., 2011), leading to a convergence rate of order  $O(\kappa/\sqrt{T})$  which is optimal in  $T$ . What is more, it can be shown that the instance of Alg. 1 in Thm. 8 achieves a rate of order  $\tilde{O}(\beta\kappa^2/T + \sigma\kappa/\sqrt{T})$  when the objective function is  $\beta$ -smooth, without requiring  $\beta$  or  $\sigma$  (the noise in the subgradients) as input. This follows by the fact that the bound in Thm. 8 is of the form  $O(\sqrt{V_T})$ —enabled by the use of FTRL-prox—and (Cutkosky, 2019, Cor. 6).

In App. C.3, we build a subroutine A (Alg. 6) based on an existing reduction due to Cutkosky (2019), such that the instance of Alg. 1 with this subroutine and  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{\mathcal{C}^\circ}$  achieves the faster rate  $\tilde{O}(\beta\kappa^2/T^2 + \sigma\kappa/\sqrt{T})$ , without knowing  $\beta$  or  $\sigma$  (see Thm. 18).

**The offline setting.** Finally, in the offline setting (a special case of the stochastic setting with zero noise, i.e.  $\sigma = 0$ ), the rates we just mentioned for general [resp.  $\beta$ -smooth] convex functions imply that at most  $\tilde{O}(d\kappa^2/\epsilon^2)$  [resp.  $\tilde{O}(d\kappa/\sqrt{\epsilon})$ ] Membership Oracle calls are required to find an  $\epsilon$ -suboptimal point. Since state-of-the-art algorithms based on the cutting plane method require  $O(d^2 \ln(1/\epsilon))$  calls to a Membership Oracle to achieve the same guarantee (Lee et al., 2018), our algorithm provides a viable alternative to the latter whenever  $d \geq \Omega(\kappa^2/\epsilon^2)$  for general convex functions or when  $d \geq \Omega(\kappa/\sqrt{\epsilon})$  for smooth functions. We note that the algorithm of Lan and Zhou (2016) has a similar guarantee to ours in the offline setting, though they require a Linear Optimization Oracle on  $\mathcal{C}$  instead of a Membership Oracle.

## 5. Discussion

In this paper, we presented a novel projection-free reduction that allowed us to turn any OCO algorithm defined on a Euclidean ball  $\mathcal{B}$ , to an algorithm on any convex set  $\mathcal{C}$  contained in  $\mathcal{B}$ , without sacrificing the performance of the original algorithm by much. Thanks to this, we were able to

build explicit algorithms that achieve optimal regret bounds in OCO without incurring expensive Euclidean projections. In particular, we swapped the latter for Gauge projections, which can be performed efficiently using a Membership Oracle; our final algorithms make at most  $O(dT \ln T)$  [resp.  $O(T \ln T)$ ] calls to such an Oracle after  $T$  rounds and guarantee a  $O(\kappa\sqrt{T})$  [resp.  $O(\kappa\sqrt{dT})$ ] regret bound, where  $\kappa := R/r$  and  $R, r > 0$  are such that  $\mathcal{B}(r) \subseteq \mathcal{C} \subseteq \mathcal{B}(R)$ .

The multiplicative factor  $\kappa$  in our regret bounds, which depends implicitly on the dimension, replaces the  $\sqrt{d}$  factor in the regret bounds of previous projection-free algorithms (see e.g. (Hazan and Kale, 2012; Hazan and Minasyan, 2020)). From Table 2, we see that  $\kappa \leq \sqrt{d}$  for many sets of interests (the bounds on  $\kappa$  are derived in App. F). Furthermore,  $\kappa$  can always be assumed to be less than  $d$  [resp.  $\sqrt{d}$ ], without loss of generality, for general [resp. centrally symmetric] convex sets by applying a certain affine reparametrization if necessary (Flaxman et al., 2005). We note that  $\kappa$  may be thought of as a measure of how well the set  $\mathcal{C}$  can be approximated by a Euclidean ball (a larger  $\kappa$  means that  $\mathcal{C}$  is not very “round”). If the set  $\mathcal{C}$  is more like a “box”, then one way to reduce  $\kappa$  is to apply our reduction (Alg. 1) with a subroutine A defined on, for example, a hypercube containing  $\mathcal{C}$  (where Euclidean projections are still cheap), rather than a Euclidean ball. We leave such investigations for future work.

For the same number of Oracle calls as in our reduction, no existing Frank-Wolfe-based algorithm guarantees a  $O(\sqrt{T})$  regret for general convex sets (unlike with our approach). However, this would be comparing apples to oranges unless we consider the computational cost of the different Oracles involved: Membership versus Linear Optimization Oracle. In Table 2, we see that the computational complexities of these Oracles are comparable for various popular sets (see also Tab. 3). In many optimization settings, the set of interest is specified via a small number of inequalities. In such cases, a Membership Oracle, which simply checks all the inequalities, typically admits a more efficient implementation than an LOO (e.g. convex-hull of permutation matrices—see Tab. 2). On the other hand, there are sets with combinatorial structure such as Matroids (where the number of inequalities defining the set can be exponential), for which performing Linear Optimization can be more efficient than Membership evaluation. In this case, Frank-Wolfe-style algorithms may be more practical. More generally, since  $(\mathcal{C}^\circ)^\circ = \mathcal{C}$  for a closed convex set  $\mathcal{C}$ , and LO on a convex set  $\mathcal{K}$  is equivalent to performing separation on its polar  $\mathcal{K}^\circ$ , it is always possible to pick a set  $\mathcal{C}$  on which LO is cheaper than performing separation/testing membership, and vice versa. For these reasons, our approach (which relies on separation/membership evaluation) may be viewed as complementary to ones that use an LOO.

**Bonus results.** Since our reduction works for any base algorithm, one can achieve different types of guarantees (e.g. a dynamic/anytime regret or one that scales with the norm of the comparator (Mhammedi and Koolen, 2020; Cutkosky, 2020)) by substituting A in Alg. 1 with an algorithm that is known to achieve the desired guarantee. We also note that our reduction can easily be extended to the setting where the losses are exp-concave instead of strongly convex using existing results due to Mhammedi et al. (2019). One can also easily extend our reduction to the Bandit setting (Chen et al., 2019; Garber and Kretzu, 2020).

<b>Domain <math>\mathcal{C}</math></b>	Upper bound on $\kappa$ post-reparam. (see App. F for param. details).	Computational Complexity of	
		LOO $\mathcal{O}_{\mathcal{C}}(\mathbf{x}; \delta)$	MO $\mathcal{M}_{\mathcal{C}}(\mathbf{x}; \delta)$
$\ell_p$ -ball in $\mathbb{R}^d$	$O(d^{\lceil 1/p-1/2 \rceil})$ , where $\mathcal{C} \subset \mathbb{R}^d$	$O(d)$	$O(d)$
Simplex $\Delta_d \in \mathbb{R}^d$	$O(d)$ , where $\mathcal{C} \subset \mathbb{R}^d$	$O(d)$	$O(d)$
Trace-norm-ball in $\mathbb{R}^{m \times n}$	$O(d^{1/4})$ , where $\mathcal{C} \subset \mathbb{R}^d$	$O(\frac{\text{nnz}(\mathbf{x})}{\sqrt{\delta}})$	Cost(SVD)
Op-norm-ball in $\mathbb{R}^{m \times n}$	$O(d^{1/4})$ , where $\mathcal{C} \subset \mathbb{R}^d$	Cost(SVD)	$O(\frac{\text{nnz}(\mathbf{x})}{\sqrt{\delta}})$
Conv-hull of Permutation Matrices in $\mathbb{R}^{n \times n}$	$O(d^{1/2})$ , where $\mathcal{C} \subset \mathbb{R}^d$	$O(n^3)$	$O(n^2)$
Convex-hull of Rotation Matrices in $\mathbb{R}^{n \times n}$	$O(d^{1/4})$ , where $\mathcal{C} \subset \mathbb{R}^d$	Cost(SVD)	Cost(SVD)
PSD matrices in $\mathbb{R}^{n \times n}$ with unit trace	$O(d)$ , where $\mathcal{C} \subset \mathbb{R}^d$	$O(\frac{\text{nnz}(\mathbf{x})}{\sqrt{\delta}})$	$O(\frac{\text{nnz}(\mathbf{x})}{\sqrt{\delta}})$
PSD matrices in $\mathbb{R}^{n \times n}$ with diagonals $\leq 1$	$O(d^{3/4})$ , where $\mathcal{C} \subset \mathbb{R}^d$	$O(\frac{\text{nnz}(\mathbf{x})}{\sqrt{\delta^2 n-3}})$	$O(\frac{\text{nnz}(\mathbf{x})}{\sqrt{\delta}})$

Table 2: If  $\mathcal{C} \subseteq \mathbb{R}^{n \times m}$ , then  $d := mn$ .  $\text{nnz}(\cdot) \equiv \#$  of non-zeros. Table entries derived in App. F.



## Acknowledgments

We thank anonymous reviewers for detailed comments that helped improve the presentation of the paper. We acknowledge support from the ONR through awards N00014-20-1-2336 and N00014-20-1-2394.

## References

- Jacob Abernethy, Chansoo Lee, Abhinav Sinha, and Ambuj Tewari. Online linear optimization via smoothing. In *Conference on Learning Theory*, pages 807–823. PMLR, 2014.
- Jacob D Abernethy, Elad Hazan, and Alexander Rakhlin. Interior-point methods for full-information and bandit online learning. *IEEE Transactions on Information Theory*, 58(7):4164–4175, 2012.
- Shalabh Bhatnagar. Adaptive newton-based multivariate smoothed functional algorithms for simulation optimization. *ACM Transactions on Modeling and Computer Simulation (TOMACS)*, 18(1): 1–35, 2007.
- Immanuel M. Bomze, Francesco Rinaldi, and Damiano Zeffiro. Frank–wolfe and friends: a journey into projection-free first-order optimization methods. *4OR*, 19(3):313–345, 2021.
- Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- Nicolo Cesa-Bianchi, Alex Conconi, and Claudio Gentile. On the generalization ability of on-line learning algorithms. *IEEE Transactions on Information Theory*, 50(9):2050–2057, 2004.
- Lin Chen, Mingrui Zhang, and Amin Karbasi. Projection-free bandit convex optimization. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2047–2056. PMLR, 2019.
- Ashok Cutkosky. Anytime online-to-batch, optimism and acceleration. In *International Conference on Machine Learning*, pages 1446–1454. PMLR, 2019.
- Ashok Cutkosky. Parameter-free, dynamic, and strongly-adaptive online learning. In *International Conference on Machine Learning*, pages 2250–2259. PMLR, 2020.
- Ashok Cutkosky and Róbert Busa-Fekete. Distributed stochastic optimization via adaptive sgd. *Advances in Neural Information Processing Systems*, 31:1910–1919, 2018.
- Ashok Cutkosky and Francesco Orabona. Black-box reductions for parameter-free online learning in banach spaces. In *Conference On Learning Theory*, pages 1493–1529. PMLR, 2018.
- John C Duchi, Peter L Bartlett, and Martin J Wainwright. Randomized smoothing for stochastic optimization. *SIAM Journal on Optimization*, 22(2):674–701, 2012.
- Rick Durrett. *Probability: theory and examples*, volume 49. Cambridge university press, 2019.
- Abraham D Flaxman, Adam Tauman Kalai, and H Brendan McMahan. Online convex optimization in the bandit setting: gradient descent without a gradient. In *Proceedings of the sixteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 385–394, 2005.

- Marguerite Frank, Philip Wolfe, et al. An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2):95–110, 1956.
- Michael P Friedlander, Ives Macedo, and Ting Kei Pong. Gauge optimization and duality. *SIAM Journal on Optimization*, 24(4):1999–2022, 2014.
- Dan Garber. Faster projection-free convex optimization over the spectrahedron. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 874–882, 2016.
- Dan Garber and Elad Hazan. A linearly convergent variant of the conditional gradient algorithm under strong convexity, with applications to online and stochastic optimization. *SIAM Journal on Optimization*, 26(3):1493–1528, 2016.
- Dan Garber and Ben Kretzu. Improved regret bounds for projection-free bandit convex optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 2196–2206. PMLR, 2020.
- JL Goffin. Affine and projective transformations in nondifferentiable optimization. *Trends in Mathematical Optimization*, pages 79–91, 1988.
- M Grötschel, L Lovász, and A Schrijver. Geometric algorithms and combinatorial optimization. *Algorithms and Combinatorics*, 1993.
- Martin Grötschel, László Lovász, and Alexander Schrijver. *Geometric algorithms and combinatorial optimization*, volume 2. Springer Science & Business Media, 2012.
- Elad Hazan. Sparse approximate solutions to semidefinite programs. In *Latin American symposium on theoretical informatics*, pages 306–316. Springer, 2008.
- Elad Hazan. Introduction to online convex optimization. *Foundations and Trends in Optimization*, 2(3-4):157–325, 2016.
- Elad Hazan. Introduction to online convex optimization. *arXiv preprint arXiv:1909.05207*, 2019.
- Elad Hazan and Satyen Kale. Projection-free online learning. In *Proceedings of the 29th International Conference on Machine Learning*, pages 1843–1850, 2012.
- Elad Hazan and Haipeng Luo. Variance-reduced and projection-free stochastic optimization. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1263–1271, New York, New York, USA, 20–22 Jun 2016. PMLR.
- Elad Hazan and Edgar Minasyan. Faster projection-free online learning. In *Conference on Learning Theory*, pages 1877–1893. PMLR, 2020.
- Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Fundamentals of convex analysis*. Springer Science & Business Media, 2004.
- Martin Jaggi. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *International Conference on Machine Learning*, pages 427–435. PMLR, 2013.

- Adam Kalai and Santosh Vempala. Efficient algorithms for online decision problems. *Journal of Computer and System Sciences*, 71(3):291–307, 2005.
- Michal Kempka, Wojciech Kotłowski, and Manfred K. Warmuth. Adaptive scale-invariant online algorithms for learning linear models. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, pages 3321–3330, 2019.
- Thomas Kerdreux. *Accelerating conditional gradient methods*. PhD thesis, Université Paris sciences et lettres, 2020.
- Wojciech Kotłowski. Scale-invariant unconstrained online learning. In *International Conference on Algorithmic Learning Theory, ALT 2017, 15-17 October 2017, Kyoto University, Kyoto, Japan*, pages 412–433, 2017.
- Ben Kretzu and Dan Garber. Revisiting projection-free online learning: the strongly convex case. In *International Conference on Artificial Intelligence and Statistics*, pages 3592–3600. PMLR, 2021.
- Jacek Kuczynski and Henryk Woźniakowski. Estimating the largest eigenvalue by the power and lanczos algorithms with a random start. *SIAM journal on matrix analysis and applications*, 13(4): 1094–1122, 1992.
- Guanghui Lan and Yi Zhou. Conditional gradient sliding for convex optimization. *SIAM Journal on Optimization*, 26(2):1379–1409, 2016.
- Guanghui Lan, Sebastian Pokutta, Yi Zhou, and Daniel Zink. Conditional accelerated lazy stochastic gradient descent. In *International Conference on Machine Learning*, pages 1965–1974. PMLR, 2017.
- Yin Tat Lee, Aaron Sidford, and Sam Chiu-wai Wong. A faster cutting plane method and its implications for combinatorial and convex optimization. In *2015 IEEE 56th Annual Symposium on Foundations of Computer Science*, pages 1049–1065. IEEE, 2015.
- Yin Tat Lee, Aaron Sidford, and Santosh S Vempala. Efficient convex optimization with membership oracles. In *Conference On Learning Theory*, pages 1292–1294. PMLR, 2018.
- Kfir Levy and Andreas Krause. Projection free online learning over smooth sets. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1458–1466. PMLR, 2019.
- Kfir Y Levy, Alp Yurtsever, and Volkan Cevher. Online adaptive methods, universality and acceleration. *Advances in Neural Information Processing Systems*, 31:6500–6509, 2018.
- Mehrdad Mahdavi, Tianbao Yang, Rong Jin, Shenghuo Zhu, and Jinfeng Yi. Stochastic gradient descent with only one projection. *Advances in neural information processing systems*, 25:494–502, 2012.
- H Brendan McMahan. A survey of algorithms and analysis for adaptive online learning. *The Journal of Machine Learning Research*, 18(1):3117–3166, 2017.

- Karola Mészáros, Alejandro H Morales, and Jessica Striker. On flow polytopes, order polytopes, and certain faces of the alternating sign matrix polytope. *Discrete & Computational Geometry*, 62(1): 128–163, 2019.
- Zakaria Mhammedi and Wouter M. Koolen. Lipschitz and comparator-norm adaptivity in online learning. In *Proceedings of Thirty Third Conference on Learning Theory*, volume 125, pages 2858–2887. PMLR, 2020.
- Zakaria Mhammedi, Wouter M Koolen, and Tim Van Erven. Lipschitz adaptivity with multiple learning rates in online learning. In *Conference on Learning Theory*, pages 2490–2511. PMLR, 2019.
- Marco Molinaro. Curvature of feasible sets in offline and online optimization. *arXiv preprint arXiv:2002.03213*, 2020.
- Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics*, 17(2):527–566, 2017.
- Francesco Orabona and Dávid Pál. Open problem: Parameter-free and scale-free online algorithms. In *Conference on Learning Theory*, pages 1659–1664. PMLR, 2016.
- Stephane Ross, Paul Mineiro, and John Langford. Normalized online learning. In *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence*, pages 537–545, 2013.
- James Saunderson, Pablo A Parrilo, and Alan S Willsky. Semidefinite descriptions of the convex hull of rotation matrices. *SIAM Journal on Optimization*, 25(3):1314–1343, 2015.
- Alexander Schrijver. *Combinatorial optimization: polyhedra and efficiency*, volume 24. Springer Science & Business Media, 2003.
- Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and trends in Machine Learning*, 4(2):107–194, 2011.
- Nathan Srebro, Karthik Sridharan, and Ambuj Tewari. Smoothness, low noise and fast rates. In J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, volume 23. Curran Associates, Inc., 2010a.
- Nathan Srebro, Karthik Sridharan, and Ambuj Tewari. Smoothness, low-noise and fast rates. In *Proceedings of the 23rd International Conference on Neural Information Processing Systems-Volume 2*, pages 2199–2207, 2010b.
- Tim Van Erven and Wouter M. Koolen. Metagrad: Multiple learning rates in online learning. In 29 (*NIPS*), pages 3666–3674, 2016.
- Tim Van Erven, Wouter M. Koolen, and Dirk Van der Hoeven. Metagrad: Adaptation using multiple learning rates in online learning. *arXiv preprint arXiv:2102.06622*, 2021.
- Olivier Wintenberger. Optimal learning with Bernstein online aggregation. *Machine Learning*, 106(1):119–141, 2017.

Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th international conference on machine learning (icml-03)*, pages 928–936, 2003.



## Appendices

**Contents**

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Setting and Definitions</b>	<b>5</b>
<b>3</b>	<b>Preliminaries: Gauge Function Approximation</b>	<b>7</b>
<b>4</b>	<b>Efficient Projection-Free Online and Stochastic Convex Optimization</b>	<b>8</b>
4.1	Algorithm for General Online Convex Optimization . . . . .	10
4.2	Algorithms for Stochastic and Offline Optimization . . . . .	11
<b>5</b>	<b>Discussion</b>	<b>11</b>
	<b>Appendices</b>	<b>19</b>
<b>A</b>	<b>Additional Notation</b>	<b>22</b>
<b>B</b>	<b>Efficient Gauge Projections using a Membership Oracle <math>\mathcal{M}_{\mathcal{C}}</math></b>	<b>22</b>
B.1	Approximating the Gauge Function $\gamma_{\mathcal{C}}$ using $\mathcal{M}_{\mathcal{C}}$ . . . . .	23
B.2	Approximating the Subgradients of $\gamma_{\mathcal{C}}$ using $\mathcal{M}_{\mathcal{C}}$ . . . . .	24
<b>C</b>	<b>Projection-Free Online and Stochastic Optimization (Detailed)</b>	<b>26</b>
C.1	A More Efficient Algorithm for General OCO . . . . .	26
C.2	Algorithm for Strongly Convex Online Optimization . . . . .	28
C.3	Efficient Projection-Free Smooth Stochastic Optimization . . . . .	30
<b>D</b>	<b>General Regret Reduction</b>	<b>32</b>
<b>E</b>	<b>Algorithm Wrapper for Scale-Invariance</b>	<b>33</b>
<b>F</b>	<b>Applying the Projection-Free Reduction in Practice</b>	<b>35</b>
F.1	$\ell_p$ -Norm Balls . . . . .	36
F.2	Simplex $\Delta_d$ . . . . .	37
F.3	Trace and Operator Norm Balls . . . . .	38
F.4	Convex-hull of Permutation Matrices . . . . .	38
F.5	Convex-hull of Rotation Matrices . . . . .	40
F.6	PSD Matrices with Unit Trace . . . . .	41
F.7	PSD Matrices with Bounded Diagonals . . . . .	43
F.8	The Flow and Matroid Polytopes . . . . .	45
<b>G</b>	<b>Adaptive OCO Algorithms</b>	<b>45</b>
<b>H</b>	<b>Linear Optimization on <math>\mathcal{C}^\circ</math> using a Membership Oracle for <math>\mathcal{C}</math></b>	<b>47</b>
H.1	Proof of Lemma 10 (Approximate Gauge Function) . . . . .	51
H.2	Proof of Proposition 11 (Approximate LO Oracle on $\mathcal{C}^\circ$ ) . . . . .	52
H.3	Proof of Lemma 13 (Efficient Stochastic LO Oracle on $\mathcal{C}^\circ$ ) . . . . .	54

<b>I</b>	<b>Proofs of the Regret Bounds and Convergence Rates</b>	<b>54</b>
I.1	Proof of Lemma 7 (Instantaneous Regret Bound) . . . . .	54
I.2	Proof of Theorem 8 (Regret Bound in High Probability using $\text{OPT}_{\mathcal{C}^\circ}$ ) . . . . .	57
I.3	Proof of Theorem 15 (Regret Bound in Expectation using $\text{OPT}_{1d, \mathcal{C}^\circ}$ ) . . . . .	58
I.4	Proof of Theorem 16 (Regret Bound in High Probability using $\text{OPT}_{1d, \mathcal{C}^\circ}$ ) . . . . .	59
I.5	Proof of Theorem 17 (The Strongly Convex Case) . . . . .	62
I.6	Proof of Theorem 18 (The Smooth Stochastic Case) . . . . .	71
<b>J</b>	<b>Technical Lemmas</b>	<b>76</b>

## Appendix A. Additional Notation

In this section, we present some additional notation that we omitted from §2 of the main body due to space. We let  $\ln_+(u) := 0 \vee \ln(u)$ , for  $u > 0$ . For  $(\mathbf{x}_t) \subset \mathbb{R}^d$  and  $t \in \mathbb{N}$ , we write  $\mathbf{x}_{1:t} := (\mathbf{x}_1, \dots, \mathbf{x}_t) \in \mathbb{R}^{d \times t}$ . For any real function  $f: \mathcal{A} \rightarrow \mathbb{R}$ , we denote by  $\arg \max_{a \in \mathcal{A}} f(a)$  [resp.  $\arg \min_{a \in \mathcal{A}} f(a)$ ] the subset  $\mathcal{K}$  of  $\mathcal{A}$  such that  $f(b) = \sup_{a \in \mathcal{A}} f(a)$  [resp.  $f(b) = \inf_{a \in \mathcal{A}} f(a)$ ], for all  $b \in \mathcal{K}$ . For any  $p \in \mathbb{R} \cup \{\infty\}$ ,  $\mathbf{u} \in \mathbb{R}^d$ , and  $r \geq 0$  we denote by  $\mathcal{B}_p(\mathbf{u}, r)$ , the  $\|\cdot\|_p$ -ball of radius  $r$  centered at  $\mathbf{u}$ , where  $\|\cdot\|_p$  denotes the  $p$ -norm. When  $\mathbf{u} = \mathbf{0}$ , we simply write  $\mathcal{B}_p(r)$  for  $\mathcal{B}_p(\mathbf{0}, r)$ . For the special case where  $p = 2$ , we let  $\|\cdot\| := \|\cdot\|_2$  and  $\mathcal{B}(\cdot) := \mathcal{B}_2(\cdot)$ . We denote by  $\Pi_{\mathcal{B}(r)}: \mathbb{R}^d \rightarrow \mathcal{B}(r)$  the Euclidean projection operator onto the ball  $\mathcal{B}(r)$ . For any set  $\mathcal{K}$ , we denote by  $\text{conv } \mathcal{K}$  [resp.  $\text{bd } \mathcal{C}$ ] the convex hull [resp. boundary] of  $\mathcal{K}$ . We let  $\iota_{\mathcal{K}}: \mathbb{R}^d \rightarrow \{0, +\infty\}$  be the indicator function of the set  $\mathcal{K}$ , where  $\iota_{\mathcal{K}}(\mathbf{w}) = 0$  if and only if  $\mathbf{w} \in \mathcal{K}$ . Finally, for any mathematical statement  $\mathcal{E}$ , we let  $\mathbb{I}_{\mathcal{E}} = 1$  if  $\mathcal{E}$  is true, and 0 otherwise.

## Appendix B. Efficient Gauge Projections using a Membership Oracle $\mathcal{M}_{\mathcal{C}}$

In this section, we build explicit algorithms that use  $\text{MEM}_{\mathcal{C}}$  (see Definition 2) to efficiently approximate  $\gamma_{\mathcal{C}}$  and its subgradients. As a result of this (and thanks to Lemma 6), we show that Gauge projections can be performed efficiently for any bounded convex set  $\mathcal{C}$  satisfying Assumption 1 using a Membership Oracle—requiring only  $\tilde{O}(d)$  calls to the latter, where  $\tilde{O}$  hides log factors in the tolerated approximation error. Thanks to Lemma 6, Gauge projections have a very simple interpretation; the projection of a point  $\mathbf{w} \notin \mathcal{C}$  onto  $\mathcal{C}$  is the point of intersection of the ray  $\{\lambda \mathbf{w} : \lambda \geq 0\}$  and the boundary of  $\mathcal{C}$  (see Figure 1). Such a point always exists under Assumption 1. Gauge projections are all we need to ensure the iterates of our new algorithms are within  $\mathcal{C}$  thanks to the carefully designed surrogate losses in Section 4.

Our starting point is Lemma 6. To prove it, we need the concept of a normal cone:

**Definition 9 (Normal Cone)** *Let  $\mathcal{C} \subseteq \mathbb{R}^d$  be a convex set. The normal cone of  $\mathcal{C}$  at  $\mathbf{u} \in \mathcal{C}$  is the set  $\mathcal{N}_{\mathcal{C}}(\mathbf{u}) := \{\mathbf{s} \in \mathcal{C} : \langle \mathbf{s}, \mathbf{w} - \mathbf{u} \rangle \leq 0, \forall \mathbf{w} \in \mathcal{C}\}$ .*

We also recall that  $\iota_{\mathcal{C}}$  denotes the indicator function of a set  $\mathcal{C}$  (see notation in App. A). With this, we now prove Lemma 6:

**Proof of Lemma 6.** Suppose that  $\gamma_{\mathcal{C}}(\mathbf{w}) > 1$  (note that by Assumption 1, we have  $\gamma_{\mathcal{C}}(\mathbf{w}) < +\infty$ ). Note that this is equivalent to  $\mathbf{w} \notin \mathcal{C}$  by the definition of the Gauge function. We will show that  $\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) \in \arg \min_{\mathbf{u} \in \mathcal{C}} \gamma_{\mathcal{C}}(\mathbf{w} - \mathbf{u})$ , which is equivalent to showing that

$$\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) \in \arg \min_{\mathbf{u} \in \mathbb{R}^d} \{\gamma_{\mathcal{C}}(\mathbf{w} - \mathbf{u}) + \iota_{\mathcal{C}}(\mathbf{u})\} = \arg \min_{\mathbf{u} \in \mathbb{R}^d} \{\sigma_{\mathcal{C}^\circ}(\mathbf{w} - \mathbf{u}) + \iota_{\mathcal{C}}(\mathbf{u})\},$$

where the last equality follows by Lemma 42-(a). Thus,  $\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) \in \arg \min_{\mathbf{u} \in \mathcal{C}} \gamma_{\mathcal{C}}(\mathbf{w} - \mathbf{u})$  if

$$\mathbf{0} \in \partial \phi(\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w})), \quad \text{where} \quad \phi(\mathbf{u}) := \sigma_{\mathcal{C}^\circ}(\mathbf{w} - \mathbf{u}) + \iota_{\mathcal{C}}(\mathbf{u}). \quad (11)$$

Using the sub-differential chain-rule and the fact that  $\partial \iota_{\mathcal{C}}(\mathbf{u}) = \mathcal{N}_{\mathcal{C}}(\mathbf{u})$  (see e.g. (Hiriart-Urruty and Lemaréchal, 2004)), where  $\mathcal{N}_{\mathcal{C}}(\mathbf{u})$  is the normal cone at  $\mathbf{u}$  (see Definition 3), we have

$$\partial \phi(\mathbf{u}) = -\partial \sigma_{\mathcal{C}^\circ}(\mathbf{w} - \mathbf{u}) + \mathcal{N}_{\mathcal{C}}(\mathbf{u}). \quad (12)$$

Let  $\mathbf{s}_* \in \partial\sigma_{\mathcal{C}^\circ}(\mathbf{w} - \mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}))$ . We will show that  $\mathbf{s}_* \in \mathcal{N}_{\mathcal{C}}(\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}))$ , where we recall that

$$\mathcal{N}_{\mathcal{C}}(\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w})) = \{\mathbf{s} \in \mathbb{R}^d : \langle \mathbf{s}, \mathbf{y} \rangle \leq \langle \mathbf{s}, \mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) \rangle, \forall \mathbf{y} \in \mathcal{C}\}, \quad (13)$$

which will imply (11) thanks to (12). By Lemma 42-(b), we have  $\mathbf{s}_* \in \partial\sigma_{\mathcal{C}^\circ}(\mathbf{w} - \mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w})) = \partial\sigma_{\mathcal{C}^\circ}(\mathbf{w}) = \arg \max_{\mathbf{s} \in \mathcal{C}^\circ} \langle \mathbf{s}, \mathbf{w} \rangle$ , and so by Lemma 42-(a), we get  $\langle \mathbf{s}_*, \mathbf{w} \rangle = \sigma_{\mathcal{C}^\circ}(\mathbf{w}) = \gamma_{\mathcal{C}}(\mathbf{w})$ . On the other hand, by definition of the polar set  $\mathcal{C}^\circ$ , we have  $\langle \mathbf{s}_*, \mathbf{y} \rangle \leq 1$ , for all  $\mathbf{y} \in \mathcal{C}$ . Using this and the fact that  $\langle \mathbf{s}_*, \mathbf{w} \rangle = \gamma_{\mathcal{C}}(\mathbf{w})$ , we get

$$\forall \mathbf{y} \in \mathcal{C}, \quad \langle \mathbf{s}_*, \mathbf{y} \rangle \leq 1 = \langle \mathbf{s}_*, \mathbf{w} \rangle / \gamma_{\mathcal{C}}(\mathbf{w}) = \langle \mathbf{s}_*, \mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) \rangle.$$

Combining this with the definition of the normal cone  $\mathcal{N}_{\mathcal{C}}(\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}))$  in (13), we get that  $\mathbf{s}_* \in \mathcal{N}_{\mathcal{C}}(\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}))$ . This shows that  $\mathbf{s}_* \in \partial\sigma_{\mathcal{C}^\circ}(\mathbf{w} - \mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w})) \cap \mathcal{N}_{\mathcal{C}}(\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}))$ , and so  $\mathbf{0} \in \partial\phi(\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}))$  by (12). This in turn implies that

$$\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) \in \arg \min_{\mathbf{u} \in \mathcal{C}} \gamma_{\mathcal{C}}(\mathbf{w} - \mathbf{u}).$$

From this result and the definition of the proximity function  $S_{\mathcal{C}}(\mathbf{w})$ , we get

$$S_{\mathcal{C}}(\mathbf{w}) = \gamma_{\mathcal{C}}(\mathbf{w} - \mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w})) = \gamma_{\mathcal{C}}(\mathbf{w} \cdot (1 - 1/\gamma_{\mathcal{C}}(\mathbf{w}))) = (1 - 1/\gamma_{\mathcal{C}}(\mathbf{w})) \cdot \gamma_{\mathcal{C}}(\mathbf{w}) = \gamma_{\mathcal{C}}(\mathbf{w}) - 1, \quad (14)$$

where the penultimate inequality follows by the current assumption that  $\gamma_{\mathcal{C}}(\mathbf{w}) > 1$  and the positive homogeneity of Gauge functions (Lemma 42-(b)). It remains to consider the case where  $\gamma_{\mathcal{C}}(\mathbf{w}) \leq 1$ . In this case, we have  $\mathbf{w} \in \mathcal{C}$ , and so

$$S_{\mathcal{C}}(\mathbf{w}) = \inf_{\mathbf{u} \in \mathcal{C}} \gamma_{\mathcal{C}}(\mathbf{w} - \mathbf{u}) = \gamma_{\mathcal{C}}(\mathbf{w} - \mathbf{w}) = 0. \quad (15)$$

In combination with (14), (15) shows that  $S_{\mathcal{C}}(\mathbf{w}) = 0 \vee (\gamma_{\mathcal{C}}(\mathbf{w}) - 1)$ , which is a convex function (since it is the maximum of two convex functions). Finally, (14) [resp. (15)] shows that  $\partial S_{\mathcal{C}}(\mathbf{w}) = \arg \max_{\mathbf{s} \in \mathcal{C}^\circ} \langle \mathbf{s}, \mathbf{w} \rangle$  when  $\mathbf{w} \notin \mathcal{C}$  ( $\equiv \gamma_{\mathcal{C}}(\mathbf{w}) > 1$ ) [resp.  $\partial S_{\mathcal{C}}(\mathbf{w}) = \{\mathbf{0}\}$  when  $\mathbf{w} \in \mathcal{C}$ ]. ■

As mentioned earlier, Lemma 6 shows the crucial fact that  $\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) \in \Pi_{\mathcal{C}}^{\mathcal{G}}(\mathbf{w})$ , for any  $\mathbf{w} \notin \mathcal{C}$ , and so to perform approximate Gauge projections (which we need to build our efficient algorithms), it suffices to approximate  $\gamma_{\mathcal{C}}$ . We technically also need to compute approximate subgradients of  $S_{\mathcal{C}}$ , which according to Lemma 6 reduces to linear optimization on  $\mathcal{C}^\circ$ . Since we do not assume access to a Linear Optimizer Oracle over  $\mathcal{C}^\circ$ , we will estimate subgradients of  $S_{\mathcal{C}}$  using our Membership Oracle  $\text{MEM}_{\mathcal{C}}$ .

### B.1. Approximating the Gauge Function $\gamma_{\mathcal{C}}$ using $\mathcal{M}_{\mathcal{C}}$

By definition of the Gauge function, we have for any  $\mathbf{w} \in \mathbb{R}^d$

$$\gamma_{\mathcal{C}}(\mathbf{w}) = \inf\{\lambda \in \mathbb{R}_{\geq 0} \mid \mathbf{w} \in \lambda\mathcal{C}\} = 1/\sup\{\nu \in \mathbb{R}_{\geq 0} \mid \nu\mathbf{w} \in \mathcal{C}\}. \quad (16)$$

Using the Membership Oracle  $\text{MEM}_{\mathcal{C}}$ , we can approximate the largest  $\nu \geq 0$  such that  $\nu\mathbf{w} \in \mathcal{C}$  via bisection, which will lead to an approximation of  $\gamma_{\mathcal{C}}(\mathbf{w})$  by (16). This is exactly what we do in Algorithm 2. We now state the guarantee of Algorithm 2 (the proof is Appendix H.1):



---

**Algorithm 2**  $\text{GAU}_{\mathcal{C}}$ : Approximate Gauge Function using Membership Oracle and the Bisection Method.

---

**Require:** Input  $(\mathbf{w}, \delta) \in \mathcal{B}(6R/5) \times (0, 1)$  with  $R$  as in (1).  $\epsilon$ -Approximate Membership Oracle  $\text{MEM}_{\mathcal{C}}(\cdot; \epsilon)$  for  $\mathcal{C}$  and  $\epsilon > 0$  (see Definition 2).

- 1: Set  $\epsilon = \delta r / (4\kappa)^2$ . //  $\kappa := R/r$ , where  $r$  and  $R$  are as in (1)
  - 2: **if**  $\text{MEM}_{\mathcal{C}}(2\mathbf{w}; \epsilon) = 1$  or  $\|\mathbf{w}\| \leq r/2$  **then**
  - 3:     Return  $\tilde{\gamma} = 0$ .
  - 4: **end if**
  - 5: Set  $\alpha = 0$ ,  $\beta = 2$ , and  $\mu = (\alpha + \beta)/2$ .
  - 6: **while**  $\beta - \alpha > \delta / (8\kappa^2)$  **do**
  - 7:     Set  $\alpha = \mu$  if  $\text{MEM}_{\mathcal{C}}(\mu\mathbf{w}; \epsilon) = 1$ ; and  $\beta = \mu$  otherwise.
  - 8:     Set  $\mu = (\alpha + \beta)/2$ .
  - 9: **end while**
  - 10: Return  $\tilde{\gamma} = (\alpha - \delta / (8\kappa^2))^{-1}$ .
- 

**Lemma 10** ( $\text{GAU}_{\mathcal{C}}$ : **Approximate Gauge Function**) *Let  $r, R > 0$  be as in (1). For any  $\delta \in (0, 1)$  and  $\mathbf{w} \in \mathcal{B}(6R/5)$ , the output  $\tilde{\gamma} = \text{GAU}_{\mathcal{C}}(\mathbf{w}; \delta)$  of Algorithm 2 satisfies*

$$\gamma_{\mathcal{C}}(\mathbf{w}) \leq \tilde{\gamma} \leq \gamma_{\mathcal{C}}(\mathbf{w}) + \delta, \quad \text{if } \gamma_{\mathcal{C}}(\mathbf{w}) \geq 9/16 \text{ or } \tilde{\gamma} \geq 1.$$

Furthermore, Alg. 2 calls the Membership Oracle  $\text{MEM}_{\mathcal{C}}(\cdot; r\delta / (4\kappa)^2)$  at most  $\lceil \log_2((4\kappa)^2 / \delta) \rceil + 1$  times.

---

**Algorithm 3**  $\text{OPT}_{\mathcal{C}^\circ}$ : Approximate Linear Optimization Algorithm on  $\mathcal{C}^\circ$ .

---

**Require:** Input point  $\mathbf{w} \in \mathcal{B}(R)$  and  $\delta \in (0, 1/3)$ .

- 1: Set  $\epsilon = \frac{r^2 \delta^3}{10^3 d^{7/2} R^2}$ ,  $\nu_1 = \frac{r\delta}{10d}$ , and  $\nu_2 = \sqrt{\frac{\epsilon \nu_1 r}{d^{1/2}}}$ . //  $r$  and  $R$  are as in (1)
  - 2: Sample  $\mathbf{u} \in \mathcal{B}_\infty(\mathbf{w}, \nu_1)$  and  $\mathbf{z} \in \mathcal{B}_\infty(\mathbf{u}, \nu_2)$  independently and uniformly at random.
  - 3: **for**  $i = 1, 2, \dots, d$  **do**
  - 4:     Let  $\mathbf{w}'_i$  and  $\mathbf{w}_i$  be the end point of the interval  $\mathcal{B}_\infty(\mathbf{u}, \nu_2) \cap \{\mathbf{z} + \lambda \mathbf{e}_i : \lambda \in \mathbb{R}\}$ .
  - 5:     Set  $\tilde{s}_i = \frac{1}{2\nu_2} (\text{GAU}_{\mathcal{C}}(\mathbf{w}'_i; \epsilon) - \text{GAU}_{\mathcal{C}}(\mathbf{w}_i; \epsilon))$ . //  $\text{GAU}_{\mathcal{C}}$  is as in Algorithm 2
  - 6: **end for**
  - 7: Set  $\tilde{\mathbf{s}} = (\tilde{s}_1, \dots, \tilde{s}_d)^\top$  and  $\tilde{\gamma} = \text{GAU}_{\mathcal{C}}(\mathbf{w}; \delta)$ .
  - 8: Return  $(\tilde{\gamma}, \tilde{\mathbf{s}})$ .
- 

## B.2. Approximating the Subgradients of $\gamma_{\mathcal{C}}$ using $\mathcal{M}_{\mathcal{C}}$

In addition to approximating  $\gamma_{\mathcal{C}}$ , we will also need to approximate its subgradients, which would then lead to approximate subgradients of the Gauge distance function  $S_{\mathcal{C}}$  by Lemma 6. The lemma also implies that approximating subgradients of  $\gamma_{\mathcal{C}}$  essentially comes down to performing linear optimization on  $\mathcal{C}^\circ$ . Algorithm 3 ( $\text{OPT}_{\mathcal{C}^\circ}$ ), which is based on (Lee et al., 2018, Alg. 2), uses  $\text{GAU}_{\mathcal{C}}$  and a random partial difference in each coordinate to approximate the subgradients of  $\gamma_{\mathcal{C}}$ . In the next proposition, we state the precise guarantee of the algorithm. The proof of the proposition, which is in

Appendix H.2, is somewhat technical and relies heavily on existing results due to Lee et al. (2018) that we restate in Appendix H.

**Proposition 11** ( $\text{OPT}_{\mathcal{C}^\circ}$ : **Approximate LOO on  $\mathcal{C}^\circ$** ) *Let  $\kappa := R/r$  with  $r, R > 0$  as in (1). For any  $\mathbf{w} \in \mathcal{B}(R)$  and  $\delta \in (0, 1/3)$ , let  $(\tilde{\gamma}, \tilde{\mathbf{s}})$  be the output of Alg. 3 with input  $(\mathbf{w}, \delta)$ . Then,  $\|\tilde{\mathbf{s}}\|_\infty < +\infty$  almost surely, and there exists a positive random variable  $\Delta \in [0, 15^2 d^4 \kappa^3 \delta^{-2}]$  satisfying  $\mathbb{E}[\Delta] \leq \delta$ , and such that if  $\tilde{\gamma} \geq 1$ , then*

$$\begin{aligned} \gamma_{\mathcal{C}}(\mathbf{u}) &\geq \gamma_{\mathcal{C}}(\mathbf{w}) + \langle \tilde{\mathbf{s}}, \mathbf{u} - \mathbf{w} \rangle - \Delta \cdot \max(1, \|\mathbf{u}\|/R), \quad \forall \mathbf{u} \in \mathbb{R}^d; \\ \|\tilde{\mathbf{s}}\|_\infty &\leq \frac{\delta}{R} + \frac{1}{r}; \quad \|\tilde{\mathbf{s}}\| \leq \frac{\Delta}{R} + \frac{1}{r}; \quad \text{and} \quad \|\tilde{\mathbf{s}}\|^2 \leq \left(\frac{2}{r} + \frac{\delta}{R}\right) \frac{\Delta}{R} + \frac{1}{r^2}. \end{aligned} \quad (17)$$

**Remark 12** (**Complexity of  $\text{OPT}_{\mathcal{C}^\circ}$** ) *In the setting of Proposition 11, Algorithm 3 ( $\text{OPT}_{\mathcal{C}^\circ}$ ) requires  $2d \cdot (\lceil \log_2((4\kappa)^2/\varepsilon) \rceil + 1)$  calls to the Membership Oracle  $\text{MEM}_{\mathcal{C}}(\cdot; r\varepsilon/(4\kappa)^2)$ , where  $\varepsilon = r^2\delta^3/(10^3 d^{7/2} R^2)$ , and one call to  $\text{MEM}_{\mathcal{C}}(\cdot; r\varepsilon/(4\kappa)^2)$ . This follows by Lemma 10 and the fact that  $\text{OPT}_{\mathcal{C}^\circ}$  makes  $2 \cdot d$  calls to  $\text{GAU}_{\mathcal{C}}(\cdot; \varepsilon)$  and one call to  $\text{GAU}_{\mathcal{C}}(\cdot; \delta)$ .*

In light of Remark 12, Proposition 11 implies that approximating the subgradient of  $\gamma_{\mathcal{C}}$  at a point  $\mathbf{w} \in \mathcal{B}(R)$  up to some random error  $\Delta \geq 0$  satisfying  $\mathbb{E}[\Delta] \leq \delta$ , requires only  $O(d \ln(d\kappa/\delta))$  calls to the Membership Oracle  $\text{MEM}_{\mathcal{C}}$ . This means that the approximation error decreases exponentially fast with the number of calls to  $\text{MEM}_{\mathcal{C}}$ , which allows us to build our efficient projection-free algorithms in Section 4.

---

**Algorithm 4**  $\text{OPT}_{1d, \mathcal{C}^\circ}$ : ‘One-Dimensional’ Stochastic Version of  $\text{OPT}_{\mathcal{C}^\circ}$ .

---

**Require:** Input point  $\mathbf{w} \in \mathcal{B}(R)$  and  $\delta \in (0, 1/3)$ .

- 1: Set  $\varepsilon = \frac{r^2\delta^3}{10^3 d^{7/2} R^2}$ ,  $\nu_1 = \frac{r\delta}{10d}$ , and  $\nu_2 = \sqrt{\frac{\varepsilon\nu_1 r}{d^{1/2}}}$ . //  $r$  and  $R$  are as in (1)
  - 2: Sample  $I \in [d]$ ,  $\mathbf{u} \in \mathcal{B}_\infty(\mathbf{w}, \nu_1)$ , and  $\mathbf{z} \in \mathcal{B}_\infty(\mathbf{u}, \nu_2)$  independently and uniformly at random.
  - 3: Let  $\mathbf{w}'_I$  and  $\mathbf{w}_I$  be the end point of the interval  $\mathcal{B}_\infty(\mathbf{u}, \nu_2) \cap \{\mathbf{z} + \lambda \mathbf{e}_I : \lambda \in \mathbb{R}\}$ .
  - 4: Set  $\tilde{\mathbf{s}}_I = \frac{1}{2\nu_2} (\text{GAU}_{\mathcal{C}}(\mathbf{w}'_I; \varepsilon) - \text{GAU}_{\mathcal{C}}(\mathbf{w}_I; \varepsilon))$ . //  $\text{GAU}_{\mathcal{C}}$  is as in Algorithm 2
  - 5: Set  $\hat{\gamma} = \text{GAU}_{\mathcal{C}}(\mathbf{w}; \delta)$  and  $\hat{\mathbf{s}} = d\tilde{\mathbf{s}}_I \cdot \mathbf{e}_I$ .
  - 6: Return  $(\hat{\gamma}, \hat{\mathbf{s}})$ .
- 

In some settings, calling a Membership Oracle  $\Omega(d)$  times per iteration might still be too expensive. A way around this is to use a stochastic version  $\text{OPT}_{1d, \mathcal{C}^\circ}$  of  $\text{OPT}_{\mathcal{C}^\circ}$  that calls  $\text{MEM}_{\mathcal{C}}$  at most  $\tilde{O}(1)$  times and has the same output as  $\text{OPT}_{\mathcal{C}^\circ}$  in expectation. Algorithm 4,  $\text{OPT}_{1d, \mathcal{C}^\circ}$ , achieves this by randomly sampling a coordinate  $I \in [d]$  for estimating the subgradient of  $\gamma_{\mathcal{C}}$  (Line 4 of Alg. 4) and using importance weights. We now state the guarantee of Algorithm 4 (the proof is in Appendix H.3):

**Lemma 13** *Let  $\delta \in (0, 1/3)$ ,  $\mathbf{w} \in \mathcal{B}(R)$ , and  $\kappa := R/r$ , where  $r, R > 0$  are as in (1). Further, let  $(\tilde{\gamma}, \tilde{\mathbf{s}}) = \text{OPT}_{\mathcal{C}^\circ}(\mathbf{w}; \delta)$  and  $(\hat{\gamma}, \hat{\mathbf{s}}) = \text{OPT}_{1d, \mathcal{C}^\circ}(\mathbf{w}; \delta)$  (Alg. 4). Then,  $\|\hat{\mathbf{s}}\| < +\infty$  a.s.;  $\tilde{\gamma} = \hat{\gamma}$ ; and if  $\hat{\gamma} \geq 1$ , it follows that*

$$\mathbb{E}[\hat{\mathbf{s}}] = \mathbb{E}[\tilde{\mathbf{s}}], \quad \text{and} \quad \mathbb{E}[\|\hat{\mathbf{s}}\|^2] \leq d \cdot (1/r + \delta/R)^2.$$

**Remark 14 (Complexity of  $\text{OPT}_{1d, \mathcal{C}^\circ}$ )** *In the setting of Lemma 13, Algorithm 4 ( $\text{OPT}_{1d, \mathcal{C}^\circ}$ ) requires  $2 \cdot (\lceil \log_2((4\kappa)^2/\varepsilon) \rceil + 1)$  calls to the Membership Oracle  $\text{MEM}_{\mathcal{C}}(\cdot; r\varepsilon/(4\kappa)^2)$ , where  $\varepsilon = r^2\delta^3/(10^3d^{7/2}R^2)$ , and one call to  $\text{MEM}_{\mathcal{C}}(\cdot; \delta)$ . This follows by Lemma 10 and the fact that  $\text{OPT}_{1d, \mathcal{C}^\circ}$  calls  $\text{GAU}_{\mathcal{C}}(\cdot; \varepsilon)$  twice and  $\text{GAU}_{\mathcal{C}}(\cdot; \delta)$  once.*

Note that  $\text{OPT}_{1d, \mathcal{C}^\circ}$  randomly selects a single coordinate  $I$  along which to estimate the subgradient of the Gauge function  $\gamma_{\mathcal{C}}$ . Generalizing this idea to  $k \leq d$  coordinates, one can build a version of  $\text{OPT}_{1d, \mathcal{C}^\circ}$ , call it  $\text{OPT}_{kd, \mathcal{C}^\circ}$ , that samples  $i \in \llbracket d/k \rrbracket$  uniformly at random and selects coordinates  $\{ik + j - 1 : j \in [k]\} \cap [d]$  along which to estimate the subgradients of  $\gamma_{\mathcal{C}}$ . In this case,  $\text{OPT}_{kd, \mathcal{C}^\circ}$  makes  $\tilde{O}(k)$  calls to the Membership Oracle and would lead to a natural trade-off between computation (i.e. Oracle calls) and regret (see discussion in Section C.1).

## Appendix C. Projection-Free Online and Stochastic Optimization (Detailed)

In this section, we provide the details of the online and stochastic optimization algorithms that we omitted from §4 due to space. In §C.1, we consider the general OCO setting where the optimization Oracle  $\mathcal{O}_{\mathcal{C}^\circ}$  in Algorithm 1 is set to  $\text{OPT}_{1d, \mathcal{C}^\circ}$ —the efficient stochastic version of  $\text{OPT}_{\mathcal{C}^\circ}$ . In §C.2, we design a subroutine A (Alg. 5) such that Algorithm 1, instantiated with this subroutine, achieves a logarithmic regret whenever the losses are strongly convex, while maintain a  $\tilde{O}(\sqrt{T})$  regret in the worst case. In §C.3, we formally define the stochastic optimization setting and design a corresponding subroutine A which, together with Alg. 1, lead to a final algorithm that adapts to the noise and smoothness of the objective function.

### C.1. A More Efficient Algorithm for General OCO

We now state the guarantee of Alg. 1 with  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{1d, \mathcal{C}^\circ}$  (the proof in Appendix I.3):

**Theorem 15** *Let  $\delta \in (0, 1/3)$  and  $\kappa := R/r$ , with  $r$  and  $R$  as in (1). Let  $(\ell_t)$  be any adversarial sequence of convex losses on  $\mathcal{C}$  and  $(\mathbf{x}_t)$  be the iterates of Alg. 1 in response to  $(\ell_t)$ . Further, let  $\mathbf{g}_s \in \partial \ell_s(\mathbf{x}_s)$ ,  $s \geq 1$ . If Alg. 1 is run with subroutine A set to FTRL-prox (Alg. 8) with parameter  $R$ ;  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{1d, \mathcal{C}^\circ}$ ; and  $\delta_t = \delta/t^2$ ,  $t \geq 1$ , then  $(\mathbf{x}_t) \subset \mathcal{C}$  and*

$$\forall \mathbf{x} \in \mathcal{C}, \quad \sum_{t=1}^T \mathbb{E}[\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle] \leq 4R \sqrt{(1 + d \cdot (\kappa + \delta)^2) \sum_{t=1}^T \mathbb{E}[\|\mathbf{g}_t\|^2]} + 6\delta R \cdot \mathbb{E} \left[ \max_{t \in [T]} \|\mathbf{g}_t\| \right]. \quad (18)$$

The instance of Algorithm 1 in Theorem 15 invokes  $\text{OPT}_{1d, \mathcal{C}^\circ}(\cdot; \delta_t)$  at each iteration  $t$ . Thus, in light of Remark 14 in App. B, the algorithm makes at most  $O(T \ln(dT\kappa/\delta))$  calls to the Membership Oracle  $\text{MEM}_{\mathcal{C}}$  after  $T$  rounds. We also remark that our choice of tolerance sequence  $(\delta_t)$  in Theorems 8 and 15 is too conservative, requiring the Membership Oracle to be more accurate than necessary to achieve a  $O(\sqrt{T})$  regret. In fact, our choice of  $(\delta_t)$  in Theorems 8 and 15 ensures that the errors involved in the approximations of the subgradients of the surrogate losses add up to a lower order term in the regret bound; i.e. the right-most terms in (9) and (18). We can choose a larger sequence of tolerances as long as the sum of the approximation errors is of order  $O(\sqrt{T})$ . Next we derive a high probability regret bound for Algorithm 1 and show that one can pick  $\delta_t$  as large as  $O(t^{-1/2})$ .

**High Probability Regret Bound.** We will now derive a high probability regret bound for Algorithm 1 where  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{1d, \mathcal{C}^\circ}$  (Alg. 4) and A is set to FTRL-prox. For the sake of simplicity, we will assume that the losses are  $B$ -Lipschitz (see (20)). Technically, this condition is not needed to derive the result in our next theorem (up to log factors in the regret), but we make it to simplify the probabilistic argument we follow in the proof of the latter. We note that the algorithm does not require knowledge of  $B$ .

**Theorem 16** *Let  $\delta \in (0, 1/3)$ ,  $B > 0$ , and  $\kappa := R/r$ , with  $r$  and  $R$  as in (1). Let  $(\ell_t)$  be any adversarial sequence of  $B$ -Lipschitz convex losses on  $\mathcal{C}$  and  $(\mathbf{x}_t)$  be the iterates of Alg. 1 in response to  $(\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t))$ . If Alg. 1 is run with subroutine A set to FTRL-prox with parameter  $R$ ;  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{1d, \mathcal{C}^\circ}$  (Alg. 4); and  $\delta_t = \delta/t^{-1/2}$ ,  $t \geq 1$ , then  $(\mathbf{x}_t) \subset \mathcal{C}$  and for all  $\rho \in (0, 1)$ ,  $T \geq d \ln(1/\rho)$ , and  $\mathbf{x} \in \mathcal{C}$ ,*

$$\mathbb{P} \left[ \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq 8RB(\kappa + \delta) \sqrt{dT(3 + 2 \ln(1/\rho))} + 2dBR(2\kappa + 2\delta + 3\delta\sqrt{T}/\rho) \right] \geq 1 - \rho.$$

The proof of the theorem is in Appendix I.4. Once again, the instance of Algorithm 1 in Theorem 16 requires only  $O(T \ln(Td\kappa/\delta))$  calls of the Membership Oracle  $\text{MEM}_{\mathcal{C}}$  (see the discussion preceding Theorem 15).

**Optimality of the bounds and computational trade-offs.** All the regret bounds in Theorems 8-16 have an optimal dependence in  $T$ . As for the dependence in  $d$ , we see that the regret bounds of Algorithm 1 in the settings of Theorems 15 and 16 can be improved by a factor of  $O(\sqrt{d})$  when making  $\tilde{O}(d)$  calls to  $\text{MEM}_{\mathcal{C}}$  per round (as in the setting of Theorem 8) or  $\tilde{O}(1)$  calls to a Separation Oracle for  $\mathcal{C}$ , if available. We also note that  $\text{OPT}_{1d, \mathcal{C}^\circ}$  (Alg. 4), which is used in the settings of Theorems 15 and 16, randomly selects a coordinate in  $[d]$  and estimates the subgradient of  $\gamma_{\mathcal{C}}$  in that coordinate's direction. A version of this algorithm, call it  $\text{OPT}_{kd, \mathcal{C}^\circ}$ , that samples  $i \in \llbracket d/k \rrbracket$  uniformly at random and selects coordinates  $\{ik + j - 1 : j \in [k]\} \cap [d]$  to estimate the subgradients of  $\gamma_{\mathcal{C}}$  would lead to a regret bound for Algorithm 1, with  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{kd, \mathcal{C}^\circ}$ , of order  $O(\sqrt{dT/k})$ , while making  $\tilde{O}(k)$  calls to  $\text{MEM}_{\mathcal{C}}$  per round. This leads to a natural trade-off between computation (Oracle calls) and regret.

Finally, we note that there is another potential dependence in  $d$  in the regret bounds through the asphericity  $\kappa$ . In Section F, we bound this quantity for the popular settings listed on Table 1. We are able to show that  $\kappa$  is often less than  $d^{1/2}$  in many settings. Nevertheless, we note that  $\kappa$  is present in our bounds because of our pessimistic upper bounds on the norms of the subgradients of the surrogate losses; i.e.  $\|\tilde{\mathbf{g}}_t\| \leq (1 + \Delta_t + \kappa)\|\mathbf{g}_t\|$ , for all  $t \geq 1$  (see Lemma 7). In fact, we do not expect the magnitude of  $(\tilde{\mathbf{g}}_t)$  to be often much larger than that of  $(\mathbf{g}_t)$  in practice; recall that  $\tilde{\mathbf{g}}_t \neq \mathbf{g}_t$  only if the iterate  $\mathbf{w}_t$  of the subroutine A is outside  $\mathcal{C}$ .

**Implications for the stochastic and offline settings.** The results we presented so far are also relevant in the stochastic and offline settings thanks to the online-to-batch conversion technique (Cesa-Bianchi et al., 2004; Shalev-Shwartz et al., 2011). In the latter settings (which we describe in more detail in Section C.3), the losses  $(\ell_t)$  are i.i.d. and satisfy  $\mathbb{E}[\ell_t] \equiv f$  for some fixed convex function  $f: \mathcal{C} \rightarrow \mathbb{R}$ . Thus, if we let  $\text{Regret}_T(\cdot)$  be the regret of Algorithm 1 in response to  $(\ell_t)$  and

$x_* \in \arg \min_{\mathbf{x} \in \mathcal{C}} f(\mathbf{x})$ , then the average iterate  $\bar{\mathbf{x}}_T$  of Alg. 1 after  $T$  rounds satisfies

$$\mathbb{E}[f(\bar{\mathbf{x}}_T)] - \inf_{\mathbf{x} \in \mathcal{C}} f(\mathbf{x}) \stackrel{(*)}{\leq} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{x}_*)] \leq \frac{\text{Regret}_T(\mathbf{x}_*)}{T}, \quad (19)$$

where  $(*)$  follows by Jensen's inequality and the fact that  $\mathbf{x}_t$  is independent of  $\ell_t$ . Plugging the regret bounds of Alg. 1 in the settings of Theorems 15-16 into (19) leads to a  $O(\sqrt{d\kappa^2/T})$  rate in stochastic and offline optimization, which is optimal in  $T^2$ . Therefore, Alg. 1 in the latter settings requires  $\tilde{O}(d\kappa^2/\epsilon^2)$  calls to the Membership Oracle to attain an  $\epsilon$ -suboptimal point in offline optimization. Thus, whenever  $d \geq \Omega(\kappa^2/\epsilon^2)$ , our algorithm is a viable alternative to those based on the cutting plane method, which require at least  $O(d^2 \ln(1/\epsilon))$  calls to a Membership Oracle to achieve the same guarantee.

When the objective function  $f$  is  $\beta$ -smooth our adaptive bounds<sup>3</sup> in Theorems 8 and 15 together with an extension of the online-to-batch conversion analysis (Cutkosky, 2019, Corollary 6) automatically imply a rate of  $O(\beta\kappa^2/T + \sigma\kappa/\sqrt{T})$ , where  $\sigma^2$  is the variance of the stochastic noise (see Section C.3 for a precise definition). Thus, in the offline setting (i.e.  $\sigma = 0$ ), our algorithm achieves the fast  $O(\beta\kappa^2/T)$  rate without knowing the smoothness parameter  $\beta$ . In Section C.3, we show how this rate can be improved further.

## C.2. Algorithm for Strongly Convex Online Optimization

In this section, we will build a subroutine A for Algorithm 1 that will enable the latter to ensure a logarithmic regret when the losses are strongly convex, while maintaining the worst-case  $O(\sqrt{T})$  regret up to log-factors. The subroutine in question is displayed in Algorithm 5. This subroutine is based on an exiting reduction due to Cutkosky and Orabona (2018) with the following key differences; I) we perform clipping of the subgradients ( $\tilde{\mathbf{g}}_t$ ) in Alg. 5; II) we slightly change the expression of the variables ( $Z_t$ ) in Line 8 of Alg 5; and III) we use FreeGrad, a *parameter-free* OCO algorithm (Mhammedi and Koolen, 2020), as the underlying OCO subroutine. These differences will allow us to make our final algorithm scale-invariant in the sense that multiplying the losses by a positive constant does not change the outputs of the algorithm. Crucially, the algorithm does not require any prior scale information on the losses unlike many existing OCO algorithms (see e.g. (Mhammedi and Koolen, 2020)).

---

2. The rate  $O(1/\sqrt{T})$  is optimal when no further assumptions on  $f$  are made (other than convexity).

3. Adaptive in the sense that they scale with  $\sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2}$  instead of  $\sqrt{T}$ .

---

**Algorithm 5** Subroutine A for Algorithm 1 for Strongly Convex Online Optimization.

---

**Require:** Parameters  $\epsilon, R > 0$ , with  $R$  as in (1) and OCO Algorithm FreeGrad (Alg. 9). Input  $(\mathbf{x}_t, f_t)$  from the Master Algorithm 1 at each round  $t \geq 1$ , where  $\mathbf{x}_t \in \mathbb{R}^d$  and  $f_t: \mathbb{R}^d \rightarrow \mathbb{R}$ .

- 1: Initialize FreeGrad with parameters  $\epsilon, R > 0$ , and set  $\mathbf{u}_0$  to FreeGrad's first output.
  - 2: Set  $\tilde{B}_0 = \epsilon$ ;  $Z_0 = 2\epsilon^2$ ; and  $\mathbf{v}_0 = \mathbf{0}$ .
  - 3: **for**  $t = 1, 2, \dots$  **do**
  - 4: Play  $\mathbf{w}_t = \mathbf{u}_{t+1} + \mathbf{v}_t/Z_t$  and observe  $\tilde{\mathbf{g}}_t \in \partial f_t(\mathbf{w}_t)$ . //  $f_t$  is specified by Alg. 1
  - 5: Set  $\tilde{B}_t = \tilde{B}_{t-1} \vee \|\tilde{\mathbf{g}}_t\|$  and  $\hat{\mathbf{g}}_t = \tilde{\mathbf{g}}_t \cdot \tilde{B}_{t-1}/\tilde{B}_t$ .
  - 6: Send linear loss  $\mathbf{w} \mapsto \langle \hat{\mathbf{g}}_t, \mathbf{w} \rangle$  to FreeGrad as the  $t$ th loss function.
  - 7: Set  $\mathbf{u}_{t+1} \in \mathcal{B}(R)$  to FreeGrad's  $(t+1)$ th output given the history  $((\mathbf{u}_i, \mathbf{w} \rightarrow \langle \hat{\mathbf{g}}_i, \mathbf{w} \rangle))_{i \leq t}$ .
  - 8: Set  $Z_t = Z_{t-1} + \|\hat{\mathbf{g}}_t\|^2 + \tilde{B}_t^2 - \tilde{B}_{t-1}^2$  and  $\mathbf{v}_t = \mathbf{v}_{t-1} + \|\hat{\mathbf{g}}_t\|^2 \mathbf{x}_t$ .  
//  $\mathbf{x}_t$  is specified by Alg. 1
  - 9: **end for**
- 

The underlying subroutine FreeGrad is displayed in Algorithm 9 in Appendix G. The key feature of FreeGrad that allows us to adapt to strong convexity is that its regret is bounded from above by  $O(\|\mathbf{w}\|\sqrt{V_T})$ , up to log-factors in  $\|\mathbf{w}\|$  and  $V_T$ , where  $V_T = \sum_{t=1}^T \|\mathbf{g}_t\|^2$  and  $(\mathbf{g}_t)$  are the observed subgradients at the iterates of the algorithm (the precise statement of the regret bound is differed to Appendix G). Thus, similar to FTRL-prox, FreeGrad also has an adaptive regret that can be much smaller than the worst-case  $O(\sqrt{T})$ ; e.g. when the losses are smooth (Srebro et al., 2010a). Furthermore, FreeGrad's regret bound scales with the norm of the comparator  $\|\mathbf{w}\|$ , and thus becomes small for comparators close to the origin. This property will be crucial to prove a logarithmic regret for strongly convex losses (Cutkosky and Orabona, 2018).

To simplify the analysis, we will state our result for  $B$ -Lipschitz losses  $(\ell_t)$  which are those that satisfy, for all  $t \geq 1$  and  $\mathbf{g}_t \in \partial \ell_t(\mathbf{x})$ :

$$\|\mathbf{g}_t\| \leq B. \quad (20)$$

We also recall that  $(\ell_t)$  are  $\mu$ -strongly convex, for  $\mu > 0$ , if for all  $t \geq 1$ ,  $\mathbf{x} \in \mathcal{C}$  and  $\mathbf{g}_t \in \partial \ell_t(\mathbf{x})$ :

$$\forall \mathbf{y} \in \mathcal{C}, \quad \ell_t(\mathbf{y}) \geq \ell_t(\mathbf{x}) + \langle \mathbf{g}_t, \mathbf{y} - \mathbf{x} \rangle + \mu \|\mathbf{x} - \mathbf{y}\|^2/2,$$

where  $\|\cdot\|$  is the Euclidean norm. With this, we now state our main result for this subsection (the proof is in Appendix I.5):

**Theorem 17** *Let  $\mu, B > 0$ ,  $\delta \in (0, 1/3)$ , and  $\kappa := R/r$ , where  $r$  and  $R$  as in (1). Suppose that Algorithm 1 is run with  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{1d, \mathcal{C}^\circ}$  (Alg. 4);  $\delta_t = \delta/t^2$ ,  $\forall t \geq 1$ ; and sub-routine A set to Alg. 5 with parameter  $\epsilon > 0$ . Then, for any adversarial sequence of convex [resp.  $\mu$ -strongly convex]  $B$ -Lipschitz losses  $(\ell_t)$  on  $\mathcal{C}$  the iterates  $(\mathbf{x}_t)$  of Algorithm 1 satisfy  $(\mathbf{x}_t) \subset \mathcal{C}$ , and for all  $T \geq 1$  and  $\mathbf{x} \in \mathcal{C}$ , we have,*

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{x})] &\leq U_T R B \sqrt{T} + R B U_T^2 / \nu, \\ \left[ \text{resp. } \sum_{t=1}^T \mathbb{E}[\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{x})] \right] &\leq \left( \frac{R}{\nu} + \frac{B}{2\mu} \right) B U_T^2, \end{aligned} \quad (21)$$

where  $U_T = O(\nu d^{1/2} \ln(e + \frac{\kappa d R T B}{\epsilon \delta}))$ ;  $\nu := 1/(R \wedge 1) + \kappa + \delta$ ; and  $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$ ,  $\forall t \geq 1$ .



We note that the instance of Algorithm 1 in the preceding theorem automatically adapts to the strong convexity constant  $\mu > 0$  of the losses  $(\ell_t)$ , where it achieves a logarithmic regret. For general convex losses, the algorithm ensures the optimal worst-case  $O(\sqrt{T})$  regret up to log-factors.

**Computational Complexity.** The instance of Algorithm 1 in Theorem 17 makes the same number of calls to the Oracle  $\text{OPT}_{1d, \mathcal{C}^\circ}$  as in the setting of Theorems 15 and 16. Thus, this instance makes at most  $O(T \ln(d\kappa T/\delta))$  calls to the Membership Oracle  $\text{MEM}_{\mathcal{C}}$ . Further, we note that the sequence  $(\delta_t)$  may be set to  $\delta_t = \delta/t^{-1/2}$  [resp.  $\delta_t = \delta_t/t$ ], for all  $t \geq 1$ , for general convex [resp. strong-convex] functions, allowing the Membership Oracle  $\text{MEM}_{\mathcal{C}}$  to be less accurate while maintaining the same regret bound as Theorem 17 up to constant factors (see discussion after Theorem 15). Finally, we restricted the losses to be  $B$ -Lipschitz in Theorem 17 only to simplify the proofs—the algorithm need not know  $B$ .

**Scale-invariance.** Though the instance of Alg. 1 in Theorem 17 does not require a bound on the norm of the gradients (an information typically required by other OL algorithms (Orabona and Pál, 2016)), the algorithm is not scale-invariant; multiplying the sequence of losses  $(\ell_t)$  by some fixed constant changes the iterates  $(\mathbf{x}_t)$  of the algorithm. In general, this is an undesirable property for OL algorithms (Orabona and Pál, 2016).

Note also that the regret bound in Theorem 17 can technically be unbounded due to the fraction  $B/\epsilon$  in the expression of  $U_T$ . This fraction can be arbitrarily large if  $\epsilon$  (a parameter of the algorithm) is too small relative to the Lipschitz constant  $B$ . Such a problematic ratio has appeared in previous works such as (Ross et al., 2013; Wintenberger, 2017; Kotłowski, 2017; Mhammedi et al., 2019; Kempka et al., 2019). To tame such a ratio, Mhammedi et al. (2019) and Mhammedi and Koolen (2020) presented a technique for certain OCO algorithms, such as MetaGrad (Van Erven and Koolen, 2016; Van Erven et al., 2021) and FreeGrad, based on the idea of restarting these algorithms whenever the ratio between the maximum norm of the observed subgradients and the norm of the initial subgradient is too large. In Appendix E, we extend this technique and present a general reduction (Algorithm 7) that makes a large class of Online Learning algorithms (including the instance of Alg. 1 in the setting of Theorem 17) scale-invariant and gets rid of problematic ratios in their regret bounds.

**Optimality and link to stochastic optimization.** The regret bounds in Theorem 17 are optimal in  $T$ . For the dependence in  $d$  in the regret bounds, the computational trade-offs, and the link to stochastic and offline optimization, see the discussion at the end of Section C.1.

### C.3. Efficient Projection-Free Smooth Stochastic Optimization

So far, we have mainly considered the setting where  $(\ell_s)$  is a sequence convex losses that may be chosen in an adversarial fashion, and the goal was to choose a sequence of iterates  $(\mathbf{x}_t) \subset \mathcal{C}$  such that the cumulative loss  $\sum_{t=1}^T \ell_t(\mathbf{x}_t)$  is small. In this subsection, we are interested in minimizing a fixed convex differentiable function  $f : \mathcal{C} \mapsto \mathbb{R}$ , with only access to a Stochastic Gradient Oracle for the function  $f$ . Formally, we assume there exists a  $\sigma > 0$  such that for any round  $t \geq 1$  and some query point  $\mathbf{x}_t \in \mathcal{C}$ , we have access to a subgradient  $\mathbf{g}_s \in \partial \ell_s(\mathbf{x}_t)$ , where

$$\ell_s(\mathbf{x}) := f(\mathbf{x}) + \langle \mathbf{x}, \boldsymbol{\xi}_s \rangle, \tag{22}$$

and  $(\boldsymbol{\xi}_s \in \mathbb{R}^d)$  are i.i.d. random vectors satisfying

$$\mathbb{E}[\boldsymbol{\xi}_s] = \mathbf{0} \quad \text{and} \quad \mathbb{E}[\|\boldsymbol{\xi}_s\|^2] \leq \sigma^2, \quad \forall s \geq 1. \tag{23}$$

We will also assume that the function  $f$  is  $\beta$ -smooth; that is,  $f$  is differentiable<sup>4</sup> on  $\mathcal{C}$  and

$$\forall \mathbf{x}, \mathbf{y} \in \mathcal{C}, \quad f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\beta}{2} \|\mathbf{y} - \mathbf{x}\|^2. \quad (24)$$

It will be instrumental to use the following consequence of (24) (see e.g. (Srebro et al., 2010b; Levy et al., 2018; Cutkosky and Busa-Fekete, 2018)); if  $f$  is  $\beta$ -smooth, then  $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}_*)\|^2 \leq 2\beta \cdot (f(\mathbf{x}) - f(\mathbf{x}_*))$ , where  $\mathbf{x}_* \in \arg \min_{\mathbf{x} \in \mathcal{C}} f(\mathbf{x})$ . Thus, when  $\mathbf{x}_*$  is in the interior of  $\mathcal{C}$ , we have  $\nabla f(\mathbf{x}_*) = \mathbf{0}$ , and so it follows that

$$\|\nabla f(\mathbf{x})\|^2 \leq 2\beta \cdot (f(\mathbf{x}) - f(\mathbf{x}_*)). \quad (25)$$

For the sake of clarity, we now summarize the assumptions we make on the loss process in this section:

**Assumption 2** *The sequence  $(\ell_s)$  in Algorithm 6 satisfies (22) with I)  $\xi_1, \xi_2, \dots \in \mathbb{R}^d$  i.i.d. vectors as in (23); and II)  $f$  is  $\beta$ -smooth, for  $\beta > 0$ , and satisfies  $\text{int}(\mathcal{C}) \cap \arg \min_{\mathbf{x} \in \mathcal{C}} f(\mathbf{x}) \neq \emptyset$ , where  $\text{int}(\mathcal{C})$  denotes the interior of  $\mathcal{C}$ .*

Without loss of generality (by making  $R$  larger if necessary), we assume that there exists  $R' \leq R$  such that:

$$\mathcal{B}(r) \subseteq \mathcal{C} \subseteq \mathcal{B}(R') \subseteq \mathcal{B}(R)/(1 + \nu), \quad \text{where } \nu := 4\sqrt{2}(1 + \kappa), \quad \text{and } \kappa := R'/r. \quad (26)$$

We now state the main result of this section (the proof is in Appendix I.6):

**Theorem 18** *Let  $\delta \in (0, 1/3)$ , and  $r, R, R', \nu$ , and  $\kappa$  be as in (26). Further, suppose that Alg. 1 is run with  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{\mathcal{C}^\circ}$  (Alg. 3);  $\delta_t = \delta/t^3, \forall t \geq 1$ ; and sub-routine A set to Alg. 6 with input  $\varepsilon > 0$ . Then, under Assumption 2, the iterates  $(\mathbf{x}_t)$  of Alg. 6, satisfy  $(\mathbf{x}_t) \subset \mathcal{C}$ , and for all  $T \geq 1$ ,*

$$\mathbb{E} [f(\mathbf{x}_T) - f(\mathbf{x}_*)] \leq \frac{2\nu\varepsilon R' + \nu^2\beta(R')^2 U_T + 3\delta R(\ln T + 6)\sqrt{\sigma^2 + 2\beta R'}}{T^2} + \frac{2\nu R' \sigma \sqrt{U_T}}{\sqrt{T}},$$

where  $\mathbf{x}_* \in \arg \min_{\mathbf{x} \in \mathcal{C}} f(\mathbf{x})$  and  $U_T := \ln(1 + 2\varepsilon^{-2}(\sigma^2 + 2R'\beta)T^3)$ .

The proof of Theorem 18 is based on an extension of a result due to Cutkosky (2019). The instance of Algorithm 1 in Theorem 18 invokes  $\text{OPT}_{\mathcal{C}^\circ}(\cdot; \delta_t)$  at each iteration  $t$ . Thus, in light of Remark 12 in App. B, the algorithm makes at most  $O(dT \ln(dT\kappa/\delta))$  calls to the Membership Oracle  $\text{MEM}_{\mathcal{C}}$  after  $T$  rounds.

**Optimality of the rate and application to the offline setting.** The rate in Theorem 18 is optimal in  $T$  and implies the fast  $O(\beta\kappa^2/T^2)$  rate in the offline smooth setting (i.e.  $\sigma = 0$ ), which is also optimal in  $T$ . Thus, Algorithm 1 in the setting of Theorem 18 reaches an  $\varepsilon$ -sub-optimal point in offline smooth optimization after  $\tilde{O}(d\kappa/\sqrt{\varepsilon})$  calls to  $\text{MEM}_{\mathcal{C}}$ . Since state-of-the-art algorithms based on the cutting plane method require  $O(d^2 \ln(1/\varepsilon))$  calls to a Membership Oracle to reach an  $\varepsilon$ -sub-optimal point (Lee et al., 2018), our algorithm provides a viable alternative to the latter whenever  $d \geq \Omega(\kappa/\sqrt{\varepsilon})$  and the objective function is smooth. As we shall see in Section F,  $\kappa$  is less

4. To avoid boundary issues, we assume (similar to (Hiriart-Urruty and Lemaréchal, 2004, Section B.4.1)) that  $\mathcal{C}$  is contained in an open set  $\Omega$  on which  $f$  is differentiable.



than  $\sqrt{d}$  in many settings of interest. We also recall that the presence of  $\kappa$  in our bounds is due to an over conservative upper bound on the norms of the subgradients of the surrogate losses, and so we expect the rates of our algorithms to scale better with  $d$  in practice.

---

**Algorithm 6** Subroutine A for Algorithm 1 for Stochastic Convex Optimization.

---

**Require:**  $r, R, R', \nu$ , and  $\kappa$  as in (26). Input  $\varepsilon > 0$ . Input  $(\mathbf{x}_t, f_t)$  from the Master Algorithm 1 at each round  $t \geq 1$ , where  $\mathbf{x}_t \in \mathbb{R}^d$  and  $f_t: \mathbb{R}^d \rightarrow \mathbb{R}$ .

- 1: Initialize FTRL-prox with parameter  $R' > 0$  and set  $\mathbf{u}_1 = \mathbf{0}$  (i.e. FTRL-prox's first output).
  - 2: Set  $\Lambda_1 = 0$ ;  $Z_0 = \varepsilon^2$ ; and  $\mathbf{w}_1 = \mathbf{0}$ .
  - 3: **for**  $t = 1, 2, \dots$  **do**
  - 4:   Output  $\mathbf{w}_t$  and observe  $\tilde{\mathbf{g}}_t \in \partial f_t(\mathbf{w}_t)$ . //  $f_t$  is provided by Alg. 1
  - 5:   Set  $\mathbf{u}_{t+1} \in \mathcal{B}(R')$  to FTRL-prox's  $(t+1)$ <sup>th</sup> output given the history  $(\mathbf{u}_i, \mathbf{w} \mapsto i \langle \tilde{\mathbf{g}}_i, \mathbf{w} \rangle)_{i \leq t}$ .
  - 6:   Set  $Z_t = Z_{t-1} + \Lambda_t \|\mathbf{g}_t\|^2$  and  $\eta_t = \nu R' / \sqrt{Z_t}$ .
  - 7:   Set  $\Lambda_{t+1} = \Lambda_t + t + 1$  and  $\mu_{t+1} = (t+1) / \Lambda_{t+1}$ .
  - 8:   Set  $\mathbf{w}_{t+1} = (1 - \mu_{t+1})(\mathbf{x}_t - \eta_t \mathbf{g}_t) + \mu_{t+1} \mathbf{u}_{t+1}$ . //  $\mathbf{x}_t$  is provided by Alg. 1
  - 9: **end for**
- 

## Appendix D. General Regret Reduction

In this section, we state and prove a regret bound for Algorithm 1, when A is a general OCO algorithm.

**Proposition 19** *Let  $\delta \in (0, 1)$ ,  $B > 0$ ,  $T \geq 1$ , and A be the sub-routine of Algorithm 1. Suppose there exists a function  $\mathcal{R}^A: \mathcal{C} \times \cup_{t=1}^{\infty} \mathbb{R}^{d \times t} \rightarrow \mathbb{R}_{\geq 0}$  such that for any adversarial sequence of convex losses  $(f_t)$  on  $\mathcal{C}$ , the iterates  $(\mathbf{w}_t)$  of A in response to  $(f_t)$ , and the subgradients  $\nabla_t \in \partial f_t(\mathbf{w}_t)$ ,  $t \in [T]$ , satisfy*

$$\sum_{t=1}^T \langle \nabla_t, \mathbf{w}_t - \mathbf{x} \rangle \leq \mathcal{R}^A(\mathbf{x}, \nabla_{1:T}), \forall \mathbf{x} \in \mathcal{C},$$

as long as  $\|\nabla_t\| \leq (1 + \frac{\delta}{t} + \kappa)B$ , for all  $t \in [T]$ . Then, for any adversarial sequence of  $B$ -Lipschitz convex losses  $(\ell_t)$ , the iterates  $(\mathbf{x}_t)$  of Alg. 1 with tolerance sequence  $\delta_t := \frac{\delta}{t^3}$ , and the subgradients  $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$ , satisfy

$$\mathbb{P} \left[ \begin{array}{l} \forall \mathbf{x} \in \mathcal{C}, \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \mathcal{R}^A(\mathbf{x}, \tilde{\mathbf{g}}_{1:T}) + \left(4 + \frac{5\delta \ln T}{\rho}\right) RB \\ \text{and } \forall t \in [T], \|\tilde{\mathbf{g}}_t\| \leq (1 + \delta/t + \kappa)B \end{array} \right] \geq 1 - 2\rho,$$

for all  $\rho \in (0, 1)$  and  $(\tilde{\mathbf{g}}_t)$  as in Algorithm 1.

**Proof** By Lemma 7, we have, for all  $\mathbf{w} \in \mathcal{C}$  and  $t \geq 1$ ,

$$\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + (2\delta_t + \Delta_t)R\|\mathbf{g}_t\|, \quad (27)$$

where  $\Delta_t \geq 0$  is a non-negative random variable satisfying  $\mathbb{E}_{t-1}[\Delta_t] \leq \delta_t$ . Thus, by the law of total expectation and Markov's inequality, we have

$$\mathbb{P}[\Delta_t \geq t^2 \delta_t / \rho] \leq \rho / t^2.$$

Let  $\mathcal{E}_T$  be the event that  $\{\Delta_t \leq t^2 \delta_t / \rho, \forall t \in [T]\}$ . By a union bound and the fact that  $\sum_{t=1}^{\infty} 1/t^2 \leq 2$ , we get that  $\mathbb{P}[\mathcal{E}_T] \geq 1 - 2\rho$ . For the rest of this proof, we will condition on the event  $\mathcal{E}_T$ . We have, by our choice of  $(\delta_t)$ ,

$$\sum_{t=1}^T \Delta_t \leq \sum_{t=1}^T t^2 \delta_t / \rho = \sum_{t=1}^T \delta / (\rho t) \leq \frac{\delta \ln T}{\rho}. \quad (28)$$

Also, note that by Lemma 7 (and the fact that  $\mathcal{E}_T$  holds), we have,

$$\|\tilde{\mathbf{g}}_t\| \leq (1 + \Delta_t + \kappa) \|\mathbf{g}_t\| \leq (1 + \delta/t + \kappa) \|\mathbf{g}_t\| \leq (1 + \delta/t + \kappa) B, \quad (29)$$

where the last inequality follows by the fact that  $\ell_t$  is  $B$ -Lipschitz. Now, by summing (27) for  $t = 1, \dots, T$ , we obtain, for all  $\mathbf{x} \in \mathcal{C}$

$$\begin{aligned} \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle &\leq \sum_{t=1}^T \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + \sum_{t=1}^T (2\delta_t + \Delta_t) R \|\mathbf{g}_t\|, \\ &\leq \mathcal{R}^A(\mathbf{x}, \tilde{\mathbf{g}}_{1:t}) + \sum_{t=1}^T (2\delta_t + \Delta_t) R \|\mathbf{g}_t\|, \\ &\leq \mathcal{R}^A(\mathbf{x}, \tilde{\mathbf{g}}_{1:t}) + (4 + 5\delta/\rho \ln(T)) RB, \end{aligned} \quad (30)$$

where (30) follows by (29) and the assumption made about the regret of Algorithm A, and the last inequality follows by (28) and the fact that  $\ell_t$  is  $B$ -Lipschitz.  $\blacksquare$

## Appendix E. Algorithm Wrapper for Scale-Invariance

In this appendix, we extend the technique of Mhammedi et al. (2019); Mhammedi and Koolen (2020) by presenting a general reduction (Algorithm 7) that makes a large class of Online Learning algorithms scale-invariant and gets rid of problematic ratios in their regret bounds. In particular, we will consider all Online Learning algorithms whose regret bound can be expressed as a map  $\mathcal{R} : \mathcal{C} \times \cup_{t \geq 1} \mathbb{R}^{d \times t} \times \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$  such that for any comparator  $\mathbf{x}$  and sequence of observed subgradients  $(\nabla_t)$ , the regret of the algorithm after  $t \geq 1$  rounds is bounded from above by  $\mathcal{R}(\mathbf{x}, \nabla_{1:t}, L_t/\epsilon)$ , where  $\epsilon > 0$  is a parameter of the algorithm and  $L_t := \epsilon \vee \max_{t \in [T]} \|\nabla_t\|$ . We note that so far we have not restricted the class of eligible algorithms by much. We will further require the following monotonicity property, which is satisfied by most popular Online Learning algorithms:

$$\mathcal{R}(\mathbf{u}, \mathbf{g}_{1:s}, p) \leq \mathcal{R}(\mathbf{u}, \mathbf{g}_{1:t}, q), \quad \text{for all } p \leq q, s \leq t, \mathbf{u} \in \mathcal{C} \text{ and } (\mathbf{g}_t) \subset \mathbb{R}^d. \quad (31)$$

With this, we now state our reduction result for scale-invariance:

**Proposition 20** *Let A be the sub-routine of Algorithm 7. Suppose there exists  $\mathcal{R}^A$  satisfying (31) and such that for any adversarial sequence of convex losses  $(f_t)$  on  $\mathcal{C}$ , the iterates  $(\mathbf{x}_t)$  of A in response to  $(f_t)$  satisfy*

$$\mathbf{x}_1 = \mathbf{0} \quad \text{and} \quad \sum_{s=1}^t \langle \nabla_s, \mathbf{x}_s - \mathbf{x} \rangle \leq \mathcal{R}^A(\mathbf{x}, \nabla_{1:t}, L_t/\epsilon), \quad \forall t \geq 1, \forall \mathbf{x} \in \mathcal{C}, \quad (32)$$

---

**Algorithm 7** Scale-Invariant Wrapper via Restarts.
 

---

**Require:** OCO Algorithm A on  $\mathcal{C}$  taking parameter  $\epsilon > 0$ .

- 1: Initialize  $\tau = 1$ ,  $\mathbf{x}_1^\tau = \mathbf{0}$ , and  $S_0 = B_0 = 0$ .
  - 2: **for**  $t = 1, 2, \dots$  **do**
  - 3:   Play  $\mathbf{y}_t = \mathbf{x}_t^\tau$  and observe  $\mathbf{g}_t \in \partial \ell_t(\mathbf{y}_t)$ .
  - 4:   Set  $B_t = B_{t-1} \vee \|\mathbf{g}_t\|$  and  $S_t = S_{t-1} + \|\mathbf{g}_t\|/B_t$ . // With the convention that  $0/0 = 0$
  - 5:   **if**  $B_t = 0$  **then**
  - 6:     Set  $\mathbf{y}_{t+1} = \mathbf{y}_t$ .
  - 7:     Continue. // play the zero vector until the first non-zero subgradients
  - 8:   **else if**  $B_{t-1} = 0$  **then**
  - 9:     Set  $\tau = t$ . //  $(B_t/B_\tau \geq S_t) \equiv (B_t/B_\tau \geq \sum_{s=\tau}^t \|\mathbf{g}_s\|/B_s)$
  - 10:    Initialize A with parameter  $\epsilon = B_\tau$ . //  $\tau$  is the new 'round 1'
  - 11:   **end if**
  - 12:   Send linear loss  $\mathbf{w} \mapsto \langle \mathbf{g}_t, \mathbf{w} \rangle$  to A.
  - 13:   Set  $\mathbf{x}_{t+1}^\tau$  to A's  $(t + 2 - \tau)$ th output given history  $((\mathbf{x}_i^\tau, \mathbf{w} \mapsto \langle \mathbf{g}_i, \mathbf{w} \rangle))_{\tau \leq i \leq t}$ .
  - 14: **end for**
- 

where  $\epsilon > 0$  is a parameter of A;  $\nabla_s \in \partial f_s(\mathbf{x}_s)$ , for all  $s \in [t]$ ; and  $L_t = \epsilon \vee \max_{s \in [t]} \|\nabla_s\|$ . Then, for any adversarial sequence of convex losses  $(\ell_t)$ , the iterates  $(\mathbf{y}_t)$  of Algorithm 7 in response to  $(\ell_t)$  satisfy

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{y}_t - \mathbf{x} \rangle \leq 2\mathcal{R}^A(\mathbf{x}, \mathbf{g}_{1:T}, T) + 4RB_T, \quad \forall T \geq 1, \forall \mathbf{x} \in \mathcal{C},$$

where  $\mathbf{g}_t, t \in [T]$ , is any subgradients in  $\partial \ell_t(\mathbf{y}_t)$  and  $B_T = \epsilon \vee \max_{t \in [T]} \|\mathbf{g}_t\|$ .

**Proof of Proposition 20.** We say that  $\tau = t$  is the start of an epoch of Algorithm 7 if the condition in Line 8 is satisfied. We use the convention that the ‘‘last epoch’’ starts at  $t = T + 1$ . Let  $\tau < t$  be the start of two consecutive epochs. Then, by the condition on Line 8 of Algorithm 7, we have

$$B_{t-1}/B_\tau \leq \sum_{s=1}^{t-1} \|\mathbf{g}_s\|/B_s \leq t - 1, \quad \text{and} \quad B_t/B_\tau > \sum_{s=1}^t \|\mathbf{g}_s\|/B_s. \quad (33)$$

Recall also that at round  $s = \tau$ , sub-routine A is initialized with  $\epsilon = B_\tau$ . Thus, by our assumption in (32), we have  $\mathbf{x}_\tau^\tau = \mathbf{0}$  and

$$\begin{aligned} \forall \mathbf{x} \in \mathcal{C}, \quad \sum_{s=\tau}^{t-1} \langle \mathbf{g}_s, \mathbf{y}_s - \mathbf{x} \rangle &= \langle \mathbf{y}_\tau, \mathbf{g}_\tau \rangle + \sum_{s=\tau}^{t-1} \langle \mathbf{g}_s, \mathbf{x}_s^\tau - \mathbf{x} \rangle \stackrel{(32)}{\leq} \langle \mathbf{y}_\tau, \mathbf{g}_\tau \rangle + \mathcal{R}^A(\mathbf{x}, \mathbf{g}_{\tau:t-1}, B_{t-1}/\epsilon), \\ &\leq RB_t + \mathcal{R}^A(\mathbf{x}, \mathbf{g}_{\tau:t-1}, t - 1), \\ &\leq RB_T + \mathcal{R}^A(\mathbf{x}, \mathbf{g}_{1:T}, T), \end{aligned} \quad (34)$$

where the last inequality follows by the fact that  $\mathcal{R}^A$  satisfies (31). If there are two epochs or less, summing (34) across the epochs leads to the desired result. Now suppose that there are more than two epochs, and let  $\tau$  [resp.  $t$ ] be the start of the ante-penultimate [resp. penultimate] epoch (recall

that the last epoch starts at  $T + 1$  by convention). By (34), the regret across these two epochs is bounded as

$$\sum_{s=\tau}^T \langle \mathbf{g}_s, \mathbf{y}_s - \mathbf{x} \rangle \leq 2RB_T + 2\mathcal{R}^A(\mathbf{x}, \mathbf{g}_{1:T}, T). \quad (35)$$

We will now bound the regret across the earlier epochs. We have

$$\sum_{s=1}^{\tau-1} \langle \mathbf{g}_s, \mathbf{y}_s - \mathbf{x} \rangle \leq 2R \sum_{s=1}^{\tau-1} \|\mathbf{g}_s\| \leq 2RB_\tau \sum_{s=1}^{\tau-1} \frac{\|\mathbf{g}_s\|}{B_s} \leq 2RB_\tau \sum_{s=1}^t \frac{\|\mathbf{g}_s\|}{B_s} \stackrel{(33)}{<} 2RB_t \leq 2RB_T. \quad (36)$$

Combining (35) and (36) leads to the desired result.  $\blacksquare$

## Appendix F. Applying the Projection-Free Reduction in Practice

In this section, we consider various popular settings where our algorithm may be applied. We note that Assumption 1, i.e. the condition that  $\mathcal{B}(r) \subseteq \mathcal{C} \subseteq \mathcal{B}(R)$  may not always be satisfied in these settings, but one can easily reparametrize the problem to satisfy the condition. We will derive explicit reparametrizations for the popular settings studied in (Hazan and Kale, 2012; Jaggi, 2013) (those in Table 3). But, first we will present a general reparametrization recipe.

Suppose that the available losses  $(f_t)$  are defined on a convex set  $\mathcal{K} \subset \mathbb{R}^d$  that does not necessarily satisfy (1) (e.g. if  $\mathcal{K}$  has an empty interior like the simplex). Further, suppose that we have a Membership Oracle  $\mathcal{M}_{\mathcal{K}}$  for  $\mathcal{K}$  and a subgradient Oracle for  $(f_t)$ . We will show that one can easily reparametrize the problem on a set  $\mathcal{C}$  satisfying (1) and whose Membership Oracle can easily be constructed from that of  $\mathcal{K}$ .

Let  $\mathcal{H}$  be the subspace generated by the span of  $\mathcal{K}$ ; that is,  $\mathcal{H} := \{\lambda \mathbf{x} : \lambda \in \mathbb{R} \text{ and } \mathbf{x} \in \text{conv} \mathcal{K}\}$ . Further, let  $(\mathbf{u}_i)_{1 \leq i \leq m}$ ,  $m \leq d$ , be an orthogonal basis of  $\mathcal{H}$  (this can be computed offline once for the given problem), and let  $\mathbf{c}$  be a point in the relative interior of  $\mathcal{K}$ . Then, one can work with the surrogate losses  $\ell_t : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$  defined by

$$\ell_t(\mathbf{x}) = \begin{cases} f_t(\mathbf{c} + \sum_{i=1}^m x_i \mathbf{u}_i), & \text{if } \mathbf{c} + \sum_{i=1}^m x_i \mathbf{u}_i \in \mathcal{K}; \\ +\infty, & \text{otherwise.} \end{cases}$$

The convexity of  $\ell_t$  follows immediately from that of  $f_t$ . We also note that since  $\mathbf{c}$  was chosen in the relative interior of  $\mathcal{K}$ , there exists  $r, R > 0$  such that the domain  $\mathcal{C} \subset \mathbb{R}^m$  of  $\ell_t$  satisfies  $\mathcal{B}(r) \subseteq \mathcal{C} \subseteq \mathcal{B}(R)$ ; that is, Assumption 1 is satisfied. We now show that a Membership Oracle for  $\mathcal{C}$  [resp. subgradient Oracle for  $\ell_t$ ] can easily be constructed from a Membership Oracle for  $\mathcal{K}$  [resp. subgradient Oracle for  $f_t$ ].

Starting with the Membership Oracle: for any  $\mathbf{x} \in \mathcal{C}$ , a Membership Oracle  $\mathcal{M}_{\mathcal{C}}$  for  $\mathcal{C}$  can be implemented as  $\mathcal{M}_{\mathcal{C}}(\mathbf{x}) = \mathcal{M}_{\mathcal{K}}(\mathbf{c} + \sum_{i=1}^m x_i \mathbf{u}_i)$ . Now, given a subgradient Oracle for  $f_t$ , it is also easy to get a subgradient Oracle for  $\ell_t$ ; for any  $\mathbf{x} \in \mathbb{R}^m$  such that  $\mathbf{y} := \sum_{i=1}^m x_i \mathbf{u}_i \in \mathcal{K}$ , we have  $\partial \ell_t(\mathbf{x}) = \{\sum_{i=1}^m \langle \mathbf{u}_i, \zeta \rangle \mathbf{u}_i : \zeta \in \partial f_t(\mathbf{y})\}$ . Using this method, one needs to perform  $O(md)$  arithmetic operations at each round to compute a subgradient of  $\ell_t$ . It is possible to avoid this computational overhead by, instead, using a finite difference approach for estimating subgradients of

Domain $\mathcal{C}$	Operation Required by		Computational Complexity of	
	LOO $\mathcal{O}_{\mathcal{C}}(\mathbf{x}; \cdot)$	MO $\mathcal{M}_{\mathcal{C}}(\mathbf{x}; \cdot)$	LOO $\mathcal{O}_{\mathcal{C}}(\cdot; \delta)$	MO $\mathcal{M}_{\mathcal{C}}(\cdot; \delta)$
$\ell_p$ -ball in $\mathbb{R}^d$	$\ \mathbf{x}\ _q$	$\ \mathbf{x}\ _p$	$O(d)$	$O(d)$
Simplex $\Delta_d \in \mathbb{R}^d$	$\ \mathbf{x}\ _{\infty}$	$\langle \mathbf{1}, \mathbf{x} \rangle$	$O(d)$	$O(d)$
Trace-norm-ball in $\mathbb{R}^{m \times n}$	$\ \mathbf{x}\ _{\text{op}}$	$\ \mathbf{x}\ _{\text{tr}}$	$O(\text{nnz}(\mathbf{x})/\sqrt{\delta})$	Cost(SVD)
Op-norm-ball in $\mathbb{R}^{m \times n}$	$\ \mathbf{x}\ _{\text{tr}}$	$\ \mathbf{x}\ _{\text{op}}$	Cost(SVD)	$O(\text{nnz}(\mathbf{x})/\sqrt{\delta})$
Conv-hull of Permutation Matrices in $\mathbb{R}^{n \times n}$		$e_i^\top \mathbf{x} \mathbf{1} = 1; \forall i$ $e_i^\top \mathbf{x}^\top \mathbf{1} = 1,$	$O(n^3)$	$O(n^2)$
Convex-hull of Rotation Matrices in $\mathbb{R}^{n \times n}$			Cost(SVD)	Cost(SVD)
PSD matrices in $\mathbb{R}^{n \times n}$ with unit trace	$\lambda_{\max}(\mathbf{x})$	$\lambda_{\min}(\mathbf{x})$ $\text{tr}(\mathbf{x})$	$O(\text{nnz}(\mathbf{x})/\sqrt{\delta})$	$O(\text{nnx}(\mathbf{x})/\sqrt{\delta})$
PSD matrices in $\mathbb{R}^{n \times n}$ with diagonals $\leq 1$		$\lambda_{\min}(\mathbf{x})$ $\max_{i \in [n]}(x_{ii})$	$O(\text{nnz}(\mathbf{x})\sqrt{n^3/\delta^5})$	$O(\text{nnx}(\mathbf{x})/\sqrt{\delta})$
The flow polytope with (#nodes, #edges)=( $d, m$ )			$\tilde{O}(d+m)$	$O(d+m)$
The Matroid polytope for Matroid $M$ ; #elem.= $d$			$\tilde{O}(d\text{Cost}(\mathcal{I}_M))$	$O(d^2\text{Cost}(\mathcal{I}_M) + d^3)$

Table 3: Computational complexity of performing linear optimization [resp. testing membership] for different sets of interest.  $\mathcal{O}_{\mathcal{C}}$  [resp.  $\mathcal{M}_{\mathcal{C}}$ ] denotes a Linear Optimization Oracle (LOO) [resp. Membership Oracle (MO)].  $\delta > 0$  represents the allowed Oracle error (see Section 2). We hide any logarithmic dependence in  $1/\delta$ . Cost(SVD) [resp. Cost( $\mathcal{I}_M$ )] represents the computational cost of performing SVD [resp. testing if a set is independent in  $M$  (see Section F)]. For  $\mathbf{x} \in \mathbb{R}^{n \times m}$ ,  $\text{nnz}(\mathbf{x})$  represents the number of non-zeros of  $\mathbf{x}$ . The details for the computational cost of the Linear Optimization Oracles listed can be found in (Hazan and Kale, 2012; Jaggi, 2013). The details for the Membership Oracle complexities can be found in Section F.

$\ell_t$ , requiring only  $2m$  calls to a value Oracle for  $f_t$  per round. Lemma 29 provides the means for doing this<sup>5</sup>.

We now show that in many popular settings there are natural parametrizations that do not require any expensive pre-processing step such as identifying a basis for the span of  $\mathcal{H}$ . In Table 4, we summarize the upper bounds we derive on the asphericity  $\kappa$  for different sets of interest after reparametrization. In Table 3, we summarize the computational complexity of a Membership Oracle for these sets.

### F.1. $\ell_p$ -Norm Balls

Consider the setting where the losses are defined on  $\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\|_p \leq 1\}$ . In this case, the span of  $\mathcal{C}$  is  $\mathbb{R}^d$  and we can pick  $\mathbf{c} = \mathbf{0}$ . For  $1 \leq p \leq 2$  Assumption 1 is satisfied with  $r = d^{1/2-1/p}$  and  $R = 1$ ; this follows by the fact that the  $\ell_p$  norm satisfies  $\|\mathbf{x}\| \leq \|\mathbf{x}\|_p \leq d^{1/p-1/2}\|\mathbf{x}\|, \forall \mathbf{x} \in \mathbb{R}^d$ . When  $2 \leq p$  Assumption 1 is satisfied with  $r = 1$  and  $R = d^{1/2-1/p}$ . Thus, in either case, the asphericity is

$$\kappa = R/r \leq d^{|1/p-1/2|},$$

5. Technically, Lemma 29 provides a way of approximating the subgradients of a function whose domain is unconstrained. However, one can extend the result to the constrained case by a careful treatment of the region around the boundary of the set.

Domain	Upper bound on $\kappa$ post-reparametrization	
$\ell_p$ -ball in $\mathbb{R}^d$	$d^{\lceil 1/p-1/2 \rceil}$	$O(d^{\lceil 1/p-1/2 \rceil})$ , where $\mathcal{C} \subset \mathbb{R}^d$
Simplex $\Delta_d \in \mathbb{R}^d$	$2d$	$O(d)$ , where $\mathcal{C} \subset \mathbb{R}^d$
Trace-norm-ball in $\mathbb{R}^{m \times n}$	$\sqrt{m \wedge n}$	$O(d^{1/4})$ , where $\mathcal{C} \subset \mathbb{R}^d$
Op-norm-ball in $\mathbb{R}^{m \times n}$	$\sqrt{m \wedge n}$	$O(d^{1/4})$ , where $\mathcal{C} \subset \mathbb{R}^d$
Conv-hull of Permutation Matrices in $\mathbb{R}^{n \times n}$	$\sqrt{5}n$	$O(d^{1/2})$ , where $\mathcal{C} \subset \mathbb{R}^d$
Convex-hull of Rotation Matrices in $\mathbb{R}^{n \times n}$	$2n^{1/2}$	$O(d^{1/4})$ , where $\mathcal{C} \subset \mathbb{R}^d$
PSD matrices in $\mathbb{R}^{n \times n}$ with unit trace	$8n^2$	$O(d)$ , where $\mathcal{C} \subset \mathbb{R}^d$
PSD matrices in $\mathbb{R}^{n \times n}$ with diagonals $\leq 1$	$4n^{3/2}$	$O(d^{3/4})$ , where $\mathcal{C} \subset \mathbb{R}^d$
The flow polytope with (#nodes, #edges)=( $d, m$ )	<i>Problem dependent</i>	
The matroid polytope for matroid $M$ ; #elem.= $d$	<i>Problem dependent</i>	

Table 4: Upper bounds on the asphericity  $\kappa$  for the different settings of Table 3 after reparametrization.

and there is no need to reparametrize. The operation needed for the Membership Oracle is simply computing  $\|\mathbf{x}\|_p$ , whereas linear optimization on  $\mathcal{C}^\circ$  amounts to evaluation the dual norm  $\|\mathbf{x}\|_q$ , where  $1/q + 1/p = 1$  (Jaggi, 2013).

## F.2. Simplex $\Delta_d$

Let  $d > 1$  and consider the setting where the losses ( $f_t$ ) are defined on the simplex  $\Delta_d := \{\mathbf{x} \in \mathbb{R}_{\geq 0}^d : \mathbf{1}^\top \mathbf{x} = 1\}$ . Since  $\Delta_d$  has an empty interior, it does not satisfy Assumption 1. However, we can easily reparametrize the problem to ensure Assumption 1.

Let  $\mathbf{c} := (\frac{1}{2} + \frac{1}{2d})\mathbf{e}_d + \sum_{i=1}^{d-1} \frac{1}{2d}\mathbf{e}_i$  and consider the sequence of reparametrized losses ( $\ell_t$ ) given by

$$\ell_t(\mathbf{x}) := f_t(\mathbf{c} + J^\top \mathbf{x}), \quad t \geq 1, \quad \text{where } J := [I_{d-1} \quad -\mathbf{1}] \in \mathbb{R}^{(d-1) \times d}. \quad (37)$$

For any  $t$ , the function  $\ell_t$  is convex and defined on the set

$$\mathcal{C} := \left\{ \mathbf{x} \in \mathbb{R}^{d-1} : x_i \geq -\frac{1}{2d} \text{ and } \mathbf{1}^\top \mathbf{x} \leq \frac{1}{2} + \frac{1}{2d} \right\}. \quad (38)$$

We now show that for this reparametrized setting, Assumption 1 holds with  $r = 1/(2d)$  and  $R = 1$ , which implies a asphericity of  $\kappa = R/r \leq 2d$ :

**Proposition 21** *The set  $\mathcal{C}$  in (38) satisfies  $\mathcal{B}(1/(2d)) \subset \mathcal{C} \subset \mathcal{B}(1)$ .*

**Proof** Note that the vertices of the set  $\mathcal{C}$  are  $\mathbf{v}_1, \dots, \mathbf{v}_d$ , where

$$\forall i \in [d-1], \quad \mathbf{v}_i = \mathbf{e}_i - \sum_{j \in [d-1]} \frac{\mathbf{e}_j}{2d}, \quad \text{and} \quad \mathbf{v}_d = - \sum_{j \in [d-1]} \frac{\mathbf{e}_j}{2d}.$$

Since  $\mathbf{v}_i \in \mathcal{B}(1)$  for all  $i \in [d]$ , we get that  $\mathcal{C} \subseteq \mathcal{B}(1)$ . We now show that  $\mathcal{B}(1/(2d)) \subseteq \mathcal{C}$ . For this, we need to find the point  $\mathbf{u}$  in the boundary of  $\mathcal{C}$  that is closest to the origin. This point must be the

orthogonal projection of the origin onto one of the  $(d - 2)$ -dimensional faces of  $\mathcal{C}$ ; there are  $d$ -many such faces corresponding to one of the inequalities defining  $\mathcal{C}$  being satisfied with equality. Thus,  $\mathbf{u}$  must satisfy one of the following:

- $\mathbf{u} = \left( \frac{1}{d-1} - \frac{1}{2d} \right) \sum_{i \in [d-1]} \mathbf{e}_i$ .
- $\exists i \in [d - 1]$ , such that  $\mathbf{u} = -\frac{\mathbf{e}_i}{2d}$ .

In all cases, we have  $\|\mathbf{u}\| \geq 1/(2d)$ , and so this shows that  $\mathcal{B}(1/(2d)) \subseteq \mathcal{C}$ . ■

It is clear from the definition of the set  $\mathcal{C}$  that the operations required to test the membership of a point  $\mathbf{x} \in \mathbb{R}^{d-1}$  are I) computing  $\langle \mathbf{1}, \mathbf{x} \rangle$ ; and II) evaluating  $x_i$ , for  $i \in [d - 1]$ . Therefore, the corresponding computational complexity is  $O(d)$ . We now show how to build a subgradient Oracle for the reparametrized losses  $(\ell_t)$ . By the chain-rule,  $\mathbf{g}$  is a subgradient of  $\ell_t$  at  $\mathbf{x}$  if and only if  $\mathbf{g} = J\zeta$ , for  $\zeta \in \partial f_t(\mathbf{c} + J^\top \mathbf{x})$ . Here,  $J\zeta$  can be evaluated in  $O(d)$ .

### F.3. Trace and Operator Norm Balls

**Trace norm ball.** Let  $m, n \geq 1$ ,  $s := m \wedge n$ , and consider the setting where the losses are defined on the trace-norm ball  $\mathcal{C} := \{\mathbf{x} \in \mathbb{R}^{m \times n} : \sum_{i=1}^s \sigma_i(\mathbf{x}) \leq 1\}$ , where  $\sigma_1(\mathbf{x}) \geq \dots \geq \sigma_s(\mathbf{x})$  are the singular values of  $\mathbf{x}$  in non-increasing order. Implementing a Membership Oracle  $\mathcal{M}_{\mathcal{C}}$  for  $\mathcal{C}$  requires computing the sum of singular values of a matrix  $\mathbf{x}$ , and so the computational complexity of  $\mathcal{M}_{\mathcal{C}}$  is at most that of performing SVD. Since the trace-norm  $\|\cdot\|_{\text{tr}}$  satisfies

$$\frac{1}{\sqrt{s}} \|\mathbf{x}\|_{\text{tr}} \leq \sqrt{\sum_{i=1}^s \sigma_i(\mathbf{x})^2} \leq \|\mathbf{x}\|_{\text{tr}},$$

and  $\sqrt{\sum_{i=1}^s \sigma_i(\mathbf{x})^2} = \sqrt{\text{tr}(\mathbf{x}\mathbf{x}^\top)} = \|\mathbf{x}\|_{\text{F}}$  is just the Euclidean norm on  $\mathbb{R}^{m \times n}$  ( $\|\cdot\|_{\text{F}}$  denotes the Frobenius norm), we have that Assumption 1 is satisfied with  $r = 1$  and  $R = \sqrt{m \wedge n}$ , and so the asphericity is  $\kappa = \sqrt{m \wedge n}$ .

**Operator norm Ball.** Let  $m, n \geq 1$ ,  $s := m \wedge n$ , and consider the setting where the losses are defined on the operator-norm ball  $\mathcal{C} := \{\mathbf{x} \in \mathbb{R}^{m \times n} : \sigma_1(\mathbf{x}) \leq 1\}$ , where  $\sigma_1(\mathbf{x}) \geq \dots \geq \sigma_s(\mathbf{x})$  are the singular values of  $\mathbf{x}$ . The operator norm  $\|\cdot\|_{\text{op}}$  is the dual to the trace-norm  $\|\cdot\|_{\text{tr}}$ . Since  $\|\mathbf{x}\|_{\text{op}}$  is the largest singular value of  $\mathbf{x}$ , we have  $\|\mathbf{x}\|_{\text{op}} \leq \|\mathbf{x}\|_{\text{F}} \leq \sqrt{s} \|\mathbf{x}\|_{\text{op}}$ , and so Assumption 1 is satisfied with  $r = (m \wedge n)^{-1/2}$  and  $R = 1$ . Implementing a Membership Oracle for  $\mathcal{C}$  requires computing the largest singular value of a given matrix  $\mathbf{x}$ . It is possible to approximate the largest singular value up to error  $\delta$  using  $O(\text{nnz}(\mathbf{x})/\sqrt{\delta})$  arithmetic operations, where  $\text{nnz}(\mathbf{x})$  represents the number of non-zeros of  $\mathbf{x}$  (Jaggi, 2013, Proposition 8). That is, the complexity of implementing a  $\delta$ -approximate Membership Oracle  $\mathcal{M}_{\mathcal{C}}(\cdot; \delta)$  (see Definition 2) is  $O(\text{nnz}(\mathbf{x})/\sqrt{\delta})$ .

### F.4. Convex-hull of Permutation Matrices

We now consider the setting where the losses  $(f_t)$  are defined on the convex-hull of permutation matrices, also known as the Birkhoff polytope  $\mathcal{K} := \{\mathbf{x} \in \mathbb{R}_{\geq 0}^{n \times n} : \mathbf{e}_i^\top \mathbf{x} \mathbf{1} = \mathbf{e}_i^\top \mathbf{x}^\top \mathbf{1} = 1\}$ . Assumption 1 is not satisfied since the *Birkhoff polytope* has an empty interior. However, as we did in

the case of the simplex, we can easily reparametrize the problem to satisfy Assumption 1. We assume that  $n > 1$ . Before presenting the reparametrized losses, we first introduce some notation. Let

$$\mathbf{c} := \sum_{1 \leq i \leq n} \left( \frac{1}{2} + \frac{1}{2n} \right) \mathbf{e}_{in} + \sum_{1 \leq j \leq n} \left( \frac{1}{2} + \frac{1}{2n} \right) \mathbf{e}_{nj} + \sum_{1 \leq i, j \leq n} \frac{\mathbf{e}_{ij}}{2n},$$

and for any  $\mathbf{x} \in \mathbb{R}^{(n-1) \times (n-1)}$  define the matrix  $\bar{\mathbf{x}} \in \mathbb{R}^{n \times n}$  such that

$$\bar{x}_{ij} := \begin{cases} x_{ij}, & \forall i, j \in [n-1]; \\ 0, & \text{otherwise.} \end{cases} \quad (39)$$

With this, consider the sequence of reparametrized losses  $(\ell_t)$  given by

$$\begin{aligned} \ell_t(\mathbf{x}) &:= f_t(\mathbf{c} + M\bar{\mathbf{x}} + \bar{\mathbf{x}}M^\top), \quad \forall \mathbf{x} \in \mathbb{R}^{(n-1) \times (n-1)}, \\ \text{where } M &:= [J^\top \quad \mathbf{0}] \in \mathbb{R}^{n \times n} \text{ and } J \text{ as in (37) with } d = n. \end{aligned} \quad (40)$$

For any  $t$ , the function  $\ell_t$  is convex and defined on the set

$$\mathcal{C} := \left\{ \mathbf{x} \in \mathbb{R}^{(n-1) \times (n-1)} : \forall i, j \in [n-1], x_{ij} \geq -\frac{1}{2n}; \sum_{k=1}^{n-1} x_{ik} \leq \frac{1}{2} + \frac{1}{2n}; \text{ and } \sum_{k=1}^{n-1} x_{kj} \leq \frac{1}{2} + \frac{1}{2n} \right\}.$$

We now show that for this reparametrized setting, Assumption 1 holds with  $r = 1/(2n)$  and  $R = \sqrt{5}/2$ . Note that the vertices of the set  $\mathcal{C}$  are  $\mathbf{v}_1, \dots, \mathbf{v}_{(n-1)^2+1}$ , where

$$\forall i \in [(n-1)^2], \quad \mathbf{v}_i = \tilde{\mathbf{e}}_i - \sum_{j \in [(n-1)^2]} \frac{\tilde{\mathbf{e}}_j}{2n}, \quad \text{and} \quad \mathbf{v}_{(n-1)^2+1} = - \sum_{j \in [(n-1)^2]} \frac{\tilde{\mathbf{e}}_j}{2n},$$

where  $\tilde{\mathbf{e}}_i = \mathbf{e}_{pq}$  and  $p, q$  are the unique integers satisfying  $i = p + (n-1)(q-1)$  and  $p, q \in [n-1]$ . Since  $\mathbf{v}_i \in \mathcal{B}(\sqrt{5}/2)$ , for all  $i \in (n-1)^2 + 1$ , we get that  $\mathcal{C} \subseteq \mathcal{B}(\sqrt{5}/2)$ . We now show that  $\mathcal{B}(1/(2n)) \subseteq \mathcal{C}$ . For this, we need to find the point  $\mathbf{u}$  in the boundary of  $\mathcal{C}$  that is closest to the origin. This point must be the orthogonal projection of the origin onto one of the  $((n-1)^2 - 1)$ -dimensional faces of  $\mathcal{C}$ ; there are  $(n^2 - 1)$ -many such faces corresponding to one of the inequalities defining  $\mathcal{C}$  being satisfied with equality. Thus,  $\mathbf{u}$  must satisfy one of the following:

- $\exists i \in [n-1]$ , such that  $\mathbf{u} = \left( \frac{1}{n-1} - \frac{1}{2n} \right) \sum_{j \in [n-1]} \mathbf{e}_{ij}$ .
- $\exists j \in [n-1]$ , such that  $\mathbf{u} = \left( \frac{1}{n-1} - \frac{1}{2n} \right) \sum_{i \in [n-1]} \mathbf{e}_{ij}$ .
- $\exists i \in [(n-1)^2]$ , such that  $\mathbf{u} = \frac{-\tilde{\mathbf{e}}_i}{2n}$ .

In all cases, we have  $\|\mathbf{u}\| \geq 1/(2n)$ , and so this shows that  $\mathcal{B}(1/(2n)) \subseteq \mathcal{C}$ . Thus, the asphericity for this reparametrized setting is

$$\kappa = R/r \leq \sqrt{5}n.$$

It is clear from the definition of the set  $\mathcal{C}$  that the operations required to test membership for a point  $\mathbf{x} \in \mathbb{R}^{n \times n}$  are I) computing  $\mathbf{e}_i^\top \mathbf{x} \mathbf{1}$  and  $\mathbf{e}_i^\top \mathbf{x}^\top \mathbf{1}$  for  $i, j \in [n]$ ; and II) evaluating  $x_{ij}$ , for  $i, j \in [n]$ . Thus the computational complexity of testing membership in  $\mathcal{C}$  is  $O(n^2)$ . We now show how to build



a subgradient Oracle for the reparametrized losses  $(\ell_t)$ . By the chain-rule,  $\mathbf{g}$  is a subgradient of  $\ell_t$  at  $\mathbf{x}$ , if and only if, for all  $i, j \in [n-1]$ ,

$$g_{ij} = 2\zeta_{ij} - \zeta_{nj} - \zeta_{in}, \quad \text{for } \zeta \in \partial f_t(\mathbf{c} + M\bar{\mathbf{x}} + \bar{\mathbf{x}}M^\top),$$

where  $\bar{\mathbf{x}}$  and  $M$  are as in (39) and (40), respectively. Since  $M$  has  $2(n-1)$  non-zero entries,  $\mathbf{g}$  can be computed in  $O(n^2)$  time (this is linear in the dimension of  $\mathcal{C}$ ).

### E.5. Convex-hull of Rotation Matrices

We now consider the setting where the losses  $(f_t)$  are defined on the convex-hull of rotation matrices; that is,

$$\mathcal{C} := \text{conv SO}(n), \quad \text{where } \text{SO}(n) := \{\mathbf{x} \in \mathbb{R}^{n \times n} : \mathbf{x}^\top \mathbf{x} = I, \det(\mathbf{x}) = 1\}.$$

This set satisfies Assumption 1 with  $r = 1 - 2/n$  and  $R = n$  (implying a asphericity of at most  $\kappa = n^2/(n-2) = O(n)$ ), and so there is not need to reparametrize as we show next:

**Proposition 22** *Let  $n > 2$ . The convex hull  $\mathcal{C}$  of Orthogonal matrices in  $\mathbb{R}^{n \times n}$  satisfies*

$$\mathcal{B}(1/2) \subseteq \mathcal{C} \subseteq \mathcal{B}(\sqrt{n}).$$

**Proof** Let  $\text{O}(n)$  the set of orthogonal matrices in  $\mathbb{R}^{n \times n}$ . By (Saunderson et al., 2015, Proposition 4.6), we have

$$\text{conv SO}(n) = (\text{conv O}(n)) \cap ((n-2)\text{SO}^-(n)^\circ), \quad (41)$$

where  $\text{SO}^-(n) := \{\mathbf{x} \in \mathbb{R}^{n \times n} : \mathbf{x}^\top \mathbf{x} = I, \det(\mathbf{x}) = -1\}$ . It is know that  $\text{conv O}(n)$  coincides with the operator-norm ball (Saunderson et al., 2015), and so we have (see paragraph on the operator norm ball above)

$$\mathcal{B}(1) \subseteq \text{conv O}(n). \quad (42)$$

We will now show that  $\mathcal{B}(1/\sqrt{n}) \subseteq \text{SO}^-(n)^\circ$  from which we conclude that  $\mathcal{B}(1/2) \subset \mathcal{B}(1 \vee (n^{1/2} - 2n^{-1/2})) \subset \text{conv SO}(n)$  using (41), (42), and the fact that  $n > 2$ . Let  $D$  be the diagonal matrix satisfying  $D_{ii} = 1$  for all  $i \in [n-1]$  and  $D_{nn} = -1$ . It is known that  $\text{SO}^-(n) = D \cdot \text{SO}(n)$  (see e.g. (Saunderson et al., 2015)). Therefore, since  $D \in \text{O}(n)$ , we have

$$\sup_{\mathbf{x} \in \text{SO}^-(n)} \|\mathbf{x}\|_{\text{op}} = \sup_{\mathbf{x} \in \text{SO}(n)} \|\mathbf{x}\|_{\text{op}} = 1.$$

Thus, by the fact that  $\|\cdot\|_{\text{F}} \leq \sqrt{n}\|\cdot\|_{\text{op}}$ , we have  $\text{SO}^-(n) \subseteq \mathcal{B}(\sqrt{n})$ , which implies that  $\mathcal{B}(1/\sqrt{n}) \subseteq \text{SO}^-(n)^\circ$ . In fact, since  $\text{SO}^-(n) \subseteq \mathcal{B}(\sqrt{n})$ , we have  $\langle \mathbf{u}, \mathbf{x} \rangle \leq 1$ , for all  $\mathbf{u} \in \mathcal{B}(1/\sqrt{n})$  and  $\mathbf{x} \in \text{SO}^-(n)$ , and so  $\mathcal{B}(1/\sqrt{n}) \subseteq \text{SO}^-(n)^\circ$  by definition of a Polar set. Combining this with the fact that  $\mathcal{B}(1) \subseteq \text{conv O}(n)$  and (41), we get that

$$\mathcal{B}(1/2) \stackrel{n > 2}{\subset} \mathcal{B}(1 \vee (n^{1/2} - 2n^{-1/2})) \subseteq \text{conv SO}(n) = \mathcal{C}.$$

We now show that  $\mathcal{C} \subseteq \mathcal{B}(n)$ . This follows by the fact that  $\mathcal{C} = \text{conv SO}(n) \subset \text{conv O}(n)$  and that  $\text{conv O}(n) \subseteq \mathcal{B}(\sqrt{n})$  since  $\text{conv O}(n)$  is the operator norm ball.  $\blacksquare$

To assess the complexity of a Membership Oracle for  $\text{conv SO}(n)$ , we use the characterization of  $\text{SO}(n)$  in (41). Also, as argued in the proof of the previous proposition,  $\text{conv O}(n)$  coincides with the operator-norm ball (Saunderson et al., 2015). Thus, in light of (41), to test if  $\mathbf{x}$  is in  $\text{conv SO}(n)$ , it suffices to test if  $\mathbf{x}$  is in the operator norm ball and in the set  $\text{SO}^-(n)^\circ$ , simultaneously. The complexity of the former test is at most that of SVD (Saunderson et al., 2015). We now show that the complexity of testing for  $\mathbf{x} \in \text{SO}^-(n)^\circ$  is also at most that of SVD up to a constant factor. First, we note that testing membership for  $\text{SO}^-(n)^\circ$  can be performed using a single call to a Linear Optimization Oracle on  $\text{SO}^-(n)$  (by leveraging the definition of a polar set). Furthermore, LO on  $\text{SO}^-(n)$  can be done using one call to a LOO on  $\text{SO}(n)$ . The latter follows by the fact that  $\text{SO}^-(n) = D \cdot \text{SO}(n)$  (see e.g. (Saunderson et al., 2015)), where  $D$  is the diagonal matrix defined in the proof of Proposition 22, and so

$$\mathcal{O}_{\text{SO}^-(n)}(\mathbf{x}) = \sup_{\mathbf{y} \in \text{SO}^-(n)} \langle \mathbf{y}, \mathbf{x} \rangle = \sup_{\mathbf{y} \in \text{SO}(n)} \langle \mathbf{y}, D\mathbf{x} \rangle = \mathcal{O}_{\text{SO}(n)}(D\mathbf{x}).$$

Finally, since the complexity of linear optimization on  $\text{SO}^-(n)$  is at most the cost of SVD (Jaggi, 2013), we conclude, in light of (41), that the complexity of a Membership Oracle for  $\text{conv SO}(n)$  is also at most that SVD up to a constant factor.

## E.6. PSD Matrices with Unit Trace

We now consider the set PSD matrices with unit trace. This set does not satisfy Assumption 1 and so we need to reparametrize. It will be useful to introduce the operator  $U: \mathbb{R}^{n(n-1)/2} \rightarrow \mathbb{R}^{n \times n}$ , where for each  $\mathbf{z} \in \mathbb{R}^{n(n-1)/2}$ ,  $U(\mathbf{z})$  is the upper-triangular matrix whose  $i$ th column is equal to  $(z_{i(i-1)/2+1}, \dots, z_{i(i+1)/2}, 0, \dots, 0)^\top \in \mathbb{R}^n$ . Further, for any  $\mathbf{x} \in \mathbb{R}^n$ , we let  $\text{diag}(\mathbf{x})$  be the matrix whose diagonal constructed from the vector  $\mathbf{x}$ , and define

$$\Theta(\mathbf{y}, \mathbf{z}) := \text{diag}(J^\top \mathbf{y}) + U(\mathbf{z}) + U(\mathbf{z})^\top,$$

for all  $\mathbf{y} \in \mathbb{R}^{n-1}$  and  $\mathbf{z} \in \mathbb{R}^{n(n-1)/2}$ , where  $J$  is as in (37) with  $d = n$ . With this, we consider the set of reparametrized losses  $(\ell_t)$  given by

$$\ell_t(\mathbf{x}) := f_t(\mathbf{c} + \Theta(\mathbf{y}, \mathbf{z})), \quad (43)$$

$$\text{where } \mathbf{x} := (\mathbf{y}, \mathbf{z}) \in \mathbb{R}^{n-1} \times \mathbb{R}^{n(n-1)/2} \quad \text{and} \quad \mathbf{c} := \left(\frac{1}{2} + \frac{1}{2n}\right) \mathbf{e}_{nn} + \sum_{i=1}^{n-1} \frac{\mathbf{e}_{ii}}{2n}.$$

For any  $t$ , the function  $\ell_t$  is convex and defined on the set

$$\mathcal{C} := \left\{ (\mathbf{y}, \mathbf{z}) \in \mathbb{R}^{n-1} \times \mathbb{R}^{n(n-1)/2} : \mathbf{c} + \Theta(\mathbf{y}, \mathbf{z}) \succeq 0 \right\}. \quad (44)$$

Furthermore, this set satisfies Assumption 1 with  $r = n^{-3/2}/4$  and  $R = 2\sqrt{n}$ , leading to an asphericity of at most  $\kappa = 8n^2$  for the set (this is linear in the dimension of  $\mathcal{C}$ ):

**Proposition 23** *The set  $\mathcal{C}$  in (44) satisfies  $\mathcal{B}(n^{-3/2}/4) \subseteq \mathcal{C} \subseteq \mathcal{B}(2\sqrt{n})$ .*

To implement a Membership Oracle for  $\mathcal{C}$  one needs to be able to test if a matrix of the form  $\mathbf{c} + \Theta(\mathbf{y}, \mathbf{z})$  is positive definite. Since this matrix is symmetric, it suffices to check if the smallest

eigenvalue of  $\Theta(\mathbf{y}, \mathbf{z})$  is non-negative. We now present a way of approximating the smallest eigenvalue of a symmetric matrix, which will then lead to an approximate Membership Oracle for  $\mathcal{C}$ . Given a symmetric matrix  $M$  and  $\delta > 0$ , first approximate its largest singular value  $\sigma_1(M)$  up to error  $\delta/2$ . This can be done using  $\tilde{O}(\text{nnz}(M)/\sqrt{\delta})$  arithmetic operations (see (Kuczyński and Woźniakowski, 1992) and (Jaggi, 2013, Proposition 8)). Next, approximate the largest singular value  $\sigma_1(M')$  of  $M' := M - \sigma_1(M) \cdot I_n$  up to error  $\delta/2$ . This also requires  $\tilde{O}(\text{nnz}(M)/\sqrt{\delta})$  arithmetic operations. Now, since the smallest eigenvalue of  $M$  is given by  $\lambda_{\min}(M) = \sigma_1(M) - \sigma_1(M')$ , we can compute a  $\delta$ -approximate value of  $\lambda_{\min}(M)$ , and thus implement a  $\delta$ -approximate Membership Oracle, using  $\tilde{O}(\text{nnz}(M)/\sqrt{\delta})$  arithmetic operations.

We now show how to build a subgradient Oracle for the reparametrized losses  $(\ell_t)$ . By the chain-rule,  $\mathbf{g} := (\mathbf{g}_y, \mathbf{g}_z)$  is a subgradient of  $\ell_t$  at  $\mathbf{x} := (\mathbf{y}, \mathbf{z})$  if and only if

$$\mathbf{g}_y = J \text{diag}^{-1}(\zeta) \quad \text{and} \quad \mathbf{g}_z := U^{-1}(\zeta), \quad \text{for } \zeta \in \partial f_t(\mathbf{c} + \Theta(\mathbf{y}, \mathbf{z})),$$

where  $J$  is as (37) and  $U^{-1}: \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n(n-1)/2}$  [resp.  $\text{diag}^{-1}: \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^n$ ] is any operator satisfying  $U^{-1} \circ U(\mathbf{z}) = \mathbf{z}$ , for all  $\mathbf{z}$  [resp.  $\text{diag}^{-1} \circ \text{diag}(\mathbf{x}) = \mathbf{x}$ , for all  $\mathbf{x}$ ]. Thus, the subgradient Oracle for  $\ell_t$  requires only an additional  $O(n^2)$  operations (this is linear in the dimension of  $\mathcal{C}$ ). We now prove Proposition 23:

**Proof of Proposition 23.** By the fact that  $\|\text{diag}(J^\top \mathbf{y})\|_F \leq 2\sqrt{n}\|\mathbf{y}\|$  and  $\|U(\mathbf{z}) + U(\mathbf{z})^\top\|_F \leq 2\|U(\mathbf{z})\|_F$ , for all  $\mathbf{y} \in \mathbb{R}^{n-1}$  and  $\mathbf{z} \in \mathbb{R}^{n(n-1)/2}$ , we have

$$\|\Theta(\mathbf{y}, \mathbf{z})\|_F / (2\sqrt{n}) \leq \|(\mathbf{y}, \mathbf{z})\| = \|\mathbf{y}\| + \|\mathbf{z}\| \leq \|\Theta(\mathbf{y}, \mathbf{z})\|_F. \quad (45)$$

We will now bound the norm of  $\Theta(\mathbf{y}, \mathbf{z})$ . For any orthogonal matrix  $H \in \mathcal{O}(n)$ ,  $\mathbf{w} \in \mathbb{R}^{n \times n}$ , and  $\boldsymbol{\lambda} \in \mathbb{R}^{n-1}$ , define

$$\begin{aligned} \Phi_H(\mathbf{w}) &:= -e_{nn}/(2n) + H^\top \mathbf{w} H; & \Psi(\boldsymbol{\lambda}) &:= \text{diag}(e_n/2 + J^\top \boldsymbol{\lambda}), \\ \text{and } \mathcal{C}' &:= \left\{ \boldsymbol{\lambda} \in \mathbb{R}^{n-1} : \lambda_i \geq \frac{-1}{2n} \text{ and } \mathbf{1}^\top \boldsymbol{\lambda} \leq \frac{1}{2} + \frac{1}{2n} \right\}, \end{aligned}$$

where  $J$  is as in (37) with  $d = n$ . Since a similarity transformation does not change the trace, or the eigenvalues for that matter, we have  $\mathbf{c} + \Phi_H \circ \Psi(\boldsymbol{\lambda}) \in \mathcal{C}$ ,  $\forall H \in \mathcal{O}(n)$ ,  $\forall \boldsymbol{\lambda} \in \mathcal{C}'$ . In particular, this implies that

$$\bigcup_{H \in \mathcal{O}(n)} \Phi_H \circ \Psi(\mathcal{C}') \subseteq \Theta(\mathcal{C}). \quad (46)$$

We will now show that for any  $H \in \mathcal{O}(n)$  and  $\boldsymbol{\lambda}$  on the boundary of  $\mathcal{C}'$ , we have  $\|\Phi_H \circ \Psi(\boldsymbol{\lambda}) - \Phi_H \circ \Psi(\mathbf{0})\|_F \geq 1/(2n)$ . This, combined with (46), would imply that  $\mathcal{B}(1/(2n)) \subseteq \Theta(\mathcal{C})$ . Let  $H \in \mathcal{O}(n)$  and  $\boldsymbol{\lambda} \in \text{bd } \mathcal{C}'$ . Since  $\Phi_H$  is an isometry with respect to the distance induced by the operator norm, we have

$$\begin{aligned} \|\Phi_H \circ \Psi(\boldsymbol{\lambda}) - \Phi_H \circ \Psi(\mathbf{0})\|_F &\geq \|\Phi_H \circ \Psi(\boldsymbol{\lambda}) - \Phi_H \circ \Psi(\mathbf{0})\|_{\text{op}} \\ &= \|\Psi(\boldsymbol{\lambda}) - \Psi(\mathbf{0})\|_{\text{op}} = \|J^\top \boldsymbol{\lambda}\| \geq \|\boldsymbol{\lambda}\| \geq 1/(2n), \end{aligned}$$

where the last inequality follows by Proposition 21. Thus, we have that  $\mathcal{B}(1/(2n)) \subseteq \Theta(\mathcal{C})$ . Combining this with (45) implies that  $\mathcal{B}(n^{-3/2}/4) \subseteq \mathcal{C}$ .

We will now show that  $\mathcal{C} \subset \mathcal{B}(2\sqrt{n})$ . For this, we will first show that (46) holds with equality. Let  $\Theta'(\cdot, \cdot) := e_{nn}/2 + \Theta(\cdot, \cdot)$ , and observe that the constraint defining the set  $\mathcal{C}$  in (44) translates to  $I_n/(2n) + \Theta'(\cdot, \cdot) \succeq 0$ . Let  $(\mathbf{y}, \mathbf{z}) \in \mathcal{C}$ . Since  $\Theta'(\mathbf{y}, \mathbf{z})$  is a real symmetric matrix, there exists an orthogonal matrix  $H$  such that  $\Lambda := H\Theta'(\mathbf{y}, \mathbf{z})H^\top$  is a diagonal matrix. Let  $\lambda_1, \dots, \lambda_n$  be the diagonal elements of  $\Lambda$ . Since these are the eigenvalues of  $\Theta'(\mathbf{y}, \mathbf{z})$ , the fact that  $I/(2n) + \Theta'(\mathbf{y}, \mathbf{z}) \succeq 0$  implies  $\lambda_i \geq -1/(2n)$ , for  $i \in [n]$ . Furthermore, since  $\text{tr}(\Theta'(\mathbf{y}, \mathbf{z})) = 1/2$ , we have  $\sum_{i=1}^n \lambda_i = 1/2$ , and so since  $\lambda_n \geq -1/(2n)$ , it follows that  $\sum_{i=1}^{n-1} \lambda_i \leq 1/2 + 1/(2n)$ . Thus, we have  $\boldsymbol{\lambda}' := (\lambda_1, \dots, \lambda_{n-1}) \in \mathcal{C}'$  and

$$e_{nn}/2 + \Theta(\mathbf{y}, \mathbf{z}) = \Theta'(\mathbf{y}, \mathbf{z}) = H^\top \text{diag}(e_n/2 + J^\top \boldsymbol{\lambda}') H = e_{nn}/2 + \Phi_H \circ \Psi(\boldsymbol{\lambda}').$$

Therefore, we have  $\Theta(\mathcal{C}) \subseteq \bigcup_{H \in \mathcal{O}(n)} \Phi_H \circ \Psi(\mathcal{C}')$ , and so by (46), we have that  $\Theta(\mathcal{C}) = \bigcup_{H \in \mathcal{O}(n)} \Phi_H \circ \Psi(\mathcal{C}')$ . Now, for any  $H \in \mathcal{O}(n)$  and  $\boldsymbol{\lambda} \in \mathcal{C}'$ , we have

$$\begin{aligned} \|\Phi_H \circ \Psi(\boldsymbol{\lambda}) - \Phi_H \circ \Psi(\mathbf{0})\|_F &\leq \sqrt{n} \|\Phi_H \circ \Psi(\boldsymbol{\lambda}) - \Phi_H \circ \Psi(\mathbf{0})\|_{\text{op}}, \\ &= \sqrt{n} \|\Psi(\boldsymbol{\lambda}) - \Psi(\mathbf{0})\|_{\text{op}}, \\ &= \sqrt{n} \|J^\top \boldsymbol{\lambda}\| \leq \sqrt{n} \|\boldsymbol{\lambda}\| + \sqrt{n} |\langle \mathbf{1}, \boldsymbol{\lambda} \rangle| \leq 2\sqrt{n}, \end{aligned}$$

where the last inequality follows by Proposition 21, and the fact that  $\langle \mathbf{1}, \boldsymbol{\lambda} \rangle \leq 1/2 + 1/(2n) \leq 1$  by definition of  $\mathcal{C}'$ . Therefore, by (45), we have  $\mathcal{C} \subseteq \mathcal{B}(2\sqrt{n})$ .  $\blacksquare$

### E.7. PSD Matrices with Bounded Diagonals

We now consider the set of PSD matrices with bounded diagonals; that is,  $\mathcal{K} := \{\mathbf{x} \in \mathbb{R}^{n \times n} : \mathbf{x} \succeq 0, \text{ and } 0 \leq x_{ii} \leq 1, \text{ for all } i \in [n]\}$ . This set does not satisfy Assumption 1 and so we need to reparametrize. As in the case of PSD matrices with unit trace, we let  $U : \mathbb{R}^{n(n-1)/2} \rightarrow \mathbb{R}^{n \times n}$  be the operator such that for each  $\mathbf{z} \in \mathbb{R}^{n(n-1)/2}$ ,  $U(\mathbf{z})$  is the upper-triangular matrix whose  $i$ th column is equal to  $(z_{i(i-1)/2+1}, \dots, z_{i(i+1)/2}, 0, \dots, 0)^\top \in \mathbb{R}^n$ . Also, for any  $\mathbf{y} \in \mathbb{R}^n$ , we let  $\text{diag}(\mathbf{y})$  be the matrix whose diagonal constructed from the vector  $\mathbf{y}$ , and define

$$\Xi(\mathbf{y}, \mathbf{z}) := \text{diag}(\mathbf{y}) + U(\mathbf{z}) + U(\mathbf{z})^\top.$$

With this, we consider the set of reparametrized losses  $(\ell_t)$  given by

$$\ell_t(\mathbf{x}) := f_t(\mathbf{c} + \Xi(\mathbf{y}, \mathbf{z})), \quad \text{where } \mathbf{x} := (\mathbf{y}, \mathbf{z}) \quad \text{and} \quad \mathbf{c} := I_n/2.$$

For any  $t$ , the function  $\ell_t$  is convex and defined on the set

$$\mathcal{C} := \left\{ (\mathbf{y}, \mathbf{z}) \in \mathbb{R}^n \times \mathbb{R}^{n(n-1)/2} : \mathbf{c} + \Xi(\mathbf{y}, \mathbf{z}) \succeq 0, \quad y_{ii} \leq 1/2, \forall i \in [n] \right\}, \quad (47)$$

Furthermore, this set satisfies Assumption 1 with  $r = 1/4$  and  $R = n^{3/2}/2$ , and so the asphericity is  $\kappa = 2n^{3/2}$ .

**Proposition 24** *The set  $\mathcal{C}$  in (47) satisfies  $\mathcal{B}(1/4) \subseteq \mathcal{C} \subseteq \mathcal{B}(n^{3/2})$ .*

To implement a Membership Oracle for  $\mathcal{C}$  one needs to be able to test if a matrix of the form  $\mathbf{c} + \Xi(\mathbf{y}, \mathbf{z})$  is positive definite. Since this matrix is symmetric, it suffices to check that the smallest eigenvalue of  $\Xi(\mathbf{y}, \mathbf{z})$  is non-negative. This can be done in the same way as in Subsection F.6 (PSD matrices with unit trace), and so a  $\delta$ -approximate Membership Oracle for  $\mathcal{C}$  can be implemented using  $\tilde{O}(\text{nnz}(\mathbf{x})/\sqrt{\delta})$  arithmetic operations for any input  $\mathbf{x} \in \mathbb{R}^n \times \mathbb{R}^{n(n-1)/2}$  and tolerance  $\delta > 0$ .

We now show how to build a subgradient Oracle for the reparametrized losses  $(\ell_t)$ . By the chain-rule,  $\mathbf{g} := (\mathbf{g}_y, \mathbf{g}_z)$  is a subgradient of  $\ell_t$  at  $\mathbf{x} := (\mathbf{y}, \mathbf{z})$  if and only if

$$\mathbf{g}_y = \text{diag}^{-1}(\zeta) \quad \text{and} \quad \mathbf{g}_z := U^{-1}(\zeta), \quad \text{for } \zeta \in \partial f_t(\mathbf{c} + \Xi(\mathbf{y}, \mathbf{z})),$$

where  $U^{-1}: \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n(n-1)/2}$  [resp.  $\text{diag}^{-1}: \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^n$ ] is an operator satisfying  $U^{-1} \circ U(\mathbf{z}) = \mathbf{z}$ , for  $\mathbf{z}$  [resp.  $\text{diag}^{-1} \circ \text{diag}(\mathbf{x}) = \mathbf{x}$ , for all  $\mathbf{x}$ ]. Thus, the subgradient Oracle for  $\ell_t$  only requires an additional  $O(n^2)$  operations (this is linear in the dimension of  $\mathcal{C}$ ).

**Proof of Proposition 24.** By the fact that  $\|\text{diag}(\mathbf{y})\|_F = \|\mathbf{y}\|$  and  $\|U(\mathbf{z}) + U(\mathbf{z})^\top\|_F \leq 2\|U(\mathbf{z})\|_F$ , for all  $\mathbf{y} \in \mathbb{R}^n$  and  $\mathbf{z} \in \mathbb{R}^{n(n-1)/2}$ , we have

$$\|\Xi(\mathbf{y}, \mathbf{z})\|_F/2 \leq \|(\mathbf{y}, \mathbf{z})\| := \|\mathbf{y}\| + \|\mathbf{z}\| \leq \|\Xi(\mathbf{y}, \mathbf{z})\|_F. \quad (48)$$

We will now bound the norm of  $\Xi(\mathbf{y}, \mathbf{z})$ . Let  $(\mathbf{y}, \mathbf{z}) \in \mathcal{C}$ . Since  $\Xi(\mathbf{y}, \mathbf{z})$  is a real symmetric matrix, there exists an orthogonal matrix  $H$  such that  $\Lambda := H\Xi(\mathbf{y}, \mathbf{z})H^\top$  is a diagonal matrix. Let  $\lambda_1, \dots, \lambda_n$  be the diagonal elements of  $\Lambda$ . Since these are also the eigenvalues of  $\Xi(\mathbf{y}, \mathbf{z})$ , the fact that  $I_n/2 + \Xi(\mathbf{y}, \mathbf{z}) \succeq 0$  implies that  $\lambda_i \geq -1/2$ , for all  $i \in [n]$ . Furthermore, since  $y_{ii} \leq 1/2$  for all  $i \in [n]$ , we have  $\text{tr}(\Xi(\mathbf{y}, \mathbf{z})) \leq n/2$  and so  $\sum_{i=1}^n \lambda_i \leq n/2$ . That is,

$$\boldsymbol{\lambda} \in \mathcal{C}' := \left\{ \mathbf{x} \in \mathbb{R}^n : x_i \geq \frac{-1}{2} \text{ and } \mathbf{1}^\top \mathbf{x} \leq \frac{n}{2} \right\}.$$

The argument above implies that

$$\Xi(\mathcal{C}) \subseteq \bigcup_{H \in \text{O}(n)} \Phi_H(\mathcal{C}'), \quad \text{where } \Phi_H(\boldsymbol{\lambda}') = H^\top \text{diag}(\boldsymbol{\lambda}') H.$$

Using that  $\|\boldsymbol{\lambda}\| \leq n$ , for all  $\boldsymbol{\lambda} \in \mathcal{C}'$ , and the fact that multiplication by an orthogonal matrix does not change the operator norm, we have

$$n \geq \sup_{\boldsymbol{\lambda} \in \mathcal{C}'} \|\boldsymbol{\lambda}\| = \sup_{\boldsymbol{\lambda} \in \mathcal{C}', H \in \text{O}(n)} \|\Phi_H(\boldsymbol{\lambda})\|_{\text{op}} \geq \sup_{\mathbf{x} \in \mathcal{C}} \|\Xi(\mathbf{x})\|_{\text{op}} \geq \sup_{\mathbf{x} \in \mathcal{C}} \|\Xi(\mathbf{x})\|_F/\sqrt{n} \geq \sup_{\mathbf{x} \in \mathcal{C}} \|\mathbf{x}\|/\sqrt{n},$$

where the last inequality follows by (48). This implies that

$$\mathcal{C} \subseteq \mathcal{B}(n^{3/2}).$$

We now show that  $\mathcal{B}(1/4) \subseteq \mathcal{C}$ . For this, we need to evaluate the quantity  $\inf_{\mathbf{x} \in \text{bd}\mathcal{C}} \|\mathbf{x}\|$ . Let  $\mathbf{x} \in \text{bd}\mathcal{C}$ . The fact that  $\mathbf{x}$  is on the boundary of  $\mathcal{C}$  implies that at least one of the inequality constraints in the definition of  $\mathcal{C}$  must be satisfied with equality; that is, one of the following must be true:

- (a)  $\exists i \in [n]$ , such that  $x_{ii} = 1/2$ .
- (b)  $\exists \mathbf{u} \in \mathbb{R}^n$ , such that  $\|\mathbf{u}\| = 1$  and  $\Xi(\mathbf{y}, \mathbf{z})\mathbf{u} = -\mathbf{c}\mathbf{u}$ .

If (a) is true, then  $\|\mathbf{x}\| \geq 1/2$ . In case (b) holds, then by (48), we have

$$\|\mathbf{x}\| \geq \|\Xi(\mathbf{y}, \mathbf{z})\|_F/2 \geq \|\Xi(\mathbf{y}, \mathbf{z})\|_{\text{op}}/2 \geq \|\Xi(\mathbf{y}, \mathbf{z})\mathbf{u}\|/2 = \|\mathbf{c}\mathbf{u}\|/2 = 1/4.$$

Since  $\mathbf{x}$  was chosen arbitrarily on the boundary of  $\mathcal{C}$ , we have that  $\mathcal{B}(1/4) \subseteq \mathcal{C}$ , which completes the proof.  $\blacksquare$

## F.8. The Flow and Matroid Polytopes

**Flow Polytope.** For this polytopes, we do not present an explicit parametrization since it is highly dependent on the specific problem at hand. We only study the complexity of the Membership Oracles for this case. The flow polytope represents the convex hull of indicator vectors corresponding to paths in a directed acyclic graph with  $d$  nodes and  $m$  edges. This polytope can be described with  $O(m + d)$  linear inequalities (Hazan and Kale, 2012; Mészáros et al., 2019). Therefore, a Membership Oracle for this polytope can be implemented using  $O(d + m)$  arithmetic operations. Linear optimization on the flow polytope has the same complexity up to log-factor (Schrijver, 2003).

**Matroid Polytope.** Same as in the case of the flow polytope, we only comment on the computational complexity of a Membership Oracle. A Matroid Polytope is the convex hull of indicator vectors corresponding to the independent sets  $A \in I$  of a matroid  $M = (E, I)$ . The polytope can be described using  $O(2^d)$  linear inequalities where  $d = |E|$ . Thus, the naive implementation of the Membership Oracle that checks all these linear inequalities would be intractable. We will present an alternative approach to designing a  $\delta$ -approximate Membership Oracle for  $M$  that requires only  $O(d^3 + d^2 \ln(d) \text{Cost}(\mathcal{I}_M)) \ln(1/\delta)$  arithmetic operations, where  $\text{Cost}(\mathcal{I}_M)$  is the computational cost (number of arithmetic operations) of testing if a subset of  $E$  is independent (i.e. an element of  $I$ ).

Let  $\mathcal{C}$  denote the matroid polytope corresponding to our matroid  $M = (E, I)$ . As we have shown in App. B, the complexity of a Membership Oracle on  $\mathcal{C}$  is the same as linear optimization on  $\mathcal{C}^\circ$ . Furthermore, a  $\delta$ -approximate Linear Optimization Oracle for  $\mathcal{C}^\circ$  can be implemented using  $O(d^3 + d \text{Cost}(\mathcal{S}_{\mathcal{C}^\circ})) \ln(1/\delta)$  arithmetic operations, where  $\text{Cost}(\mathcal{S}_{\mathcal{C}^\circ})$  is the computational cost (number of arithmetic operations) of a Separation Oracle on  $\mathcal{C}^\circ$  (Lee et al., 2018). The fact that  $(\mathcal{C}^\circ)^\circ = \mathcal{C}$  for a closed convex set and our results from App. B (see also (Grötschel et al., 1993; Molinaro, 2020)) imply that the complexity of  $\mathcal{S}_{\mathcal{C}^\circ}$  is the same as linear optimization on  $\mathcal{C}$ . The latter can be performed using  $O(d \ln(d) \text{Cost}(\mathcal{I}_M))$  arithmetic operations (Schrijver, 2003, Section 40.1). All in all, a Membership Oracle for  $\mathcal{C}$  can be implemented using  $O(d^3 + d^2 \ln(d) \text{Cost}(\mathcal{I}_M))$  arithmetic operations (omitting log-factors in  $1/\delta$ ).

## Appendix G. Adaptive OCO Algorithms

We now present two algorithms for Online Convex Optimization that we will build on to derive our results.

---

**Algorithm 8** FTRL-proximal on  $\mathcal{B}(R)$  (McMahan, 2017, Algorithm 2 & Section 3.3)

---

**Require:** Radius  $R$ .

- 1: Initialize  $\mathbf{w}_1 = \mathbf{0}$ ,  $\mathbf{G}_0 = \mathbf{0}$ , and  $V_0 = 0$ .
  - 2: Initialize  $\eta_0 = \infty$  and  $\Sigma_0 = 0$ ,
  - 3: **for**  $t = 1, 2, \dots$  **do**
  - 4:   Play  $\mathbf{w}_t$  and observe  $\nabla_t \in \partial \ell_t(\mathbf{w}_t)$ .
  - 5:   Set  $\eta_t = \sqrt{2}R/\sqrt{V_t}$ ,  $\sigma_t = 1/\eta_t - 1/\eta_{t-1}$ , and  $\Sigma_t = \Sigma_{t-1} + \sigma_t$ .  
      // With the convention  $1/\infty := 0$
  - 6:   Set  $\mathbf{G}_t = \mathbf{G}_{t-1} + \nabla_t$  and  $V_t = V_{t-1} + \|\nabla_t\|^2$ .
  - 7:   Set  $\mathbf{w}_{t+1} = \Pi_{\mathcal{B}(R)}(\tilde{\mathbf{w}}_{t+1})$ , where  $\tilde{\mathbf{w}}_{t+1} := (-\mathbf{G}_t + \sum_{s=1}^t \sigma_s \mathbf{w}_s)/\Sigma_t$ .
  - 8:   // The vector  $\mathbf{w}_{t+1}$  above satisfies  $\mathbf{w}_{t+1} \in \arg \min_{\mathbf{w} \in \mathcal{B}(R)} \langle \mathbf{G}_t, \mathbf{w} \rangle + \sum_{s=1}^t \sigma_s^2 \|\mathbf{w} - \mathbf{w}_s\|^2/2$ .
  - 9: **end for**
- 

The first algorithm, FTRL-prox(Alg. 8), requires  $R$  as input, but does not require an upper bound on the norm of the input loss vectors ( $\nabla_t$ ) (i.e. the algorithm adapts to the norm of the loss vectors). We now state the guarantee of FTRL-prox which follows from (McMahan, 2017, Theorem 2 & Section 3.3) with the choice of learning rate  $\eta_t := \sqrt{2}R/\sqrt{V_t}$ , where  $V_t := \sum_{s=1}^t \|\nabla_s\|^2$ :

**Proposition 25 (FTRL-proximal’s regret)** *For any adversarial sequence of convex losses  $(\ell_t)$  on  $\mathcal{B}(R)$ , the iterates  $(\mathbf{w}_t)$  of FTRL-prox (Alg. 8) with parameter  $R > 0$  in response to  $(\ell_t)$ , satisfy for all  $T \geq 1$  and  $\mathbf{w} \in \mathcal{B}(R)$ ,*

$$\sum_{t=1}^T (\ell_t(\mathbf{w}_t) - \ell_t(\mathbf{w})) \leq \sum_{t=1}^T \langle \nabla_t, \mathbf{w}_t - \mathbf{w} \rangle \leq 2R\sqrt{2V_T},$$

where  $\nabla_t$ ,  $t \geq 1$ , is any subgradient in  $\partial \ell_t(\mathbf{w}_t)$  and  $V_T := \sum_{t=1}^T \|\nabla_t\|^2$ .

---

**Algorithm 9** FreeGrad (Mhammedi and Koolen, 2020) with the unconstrained-to-constrained reduction due to Cutkosky (2020).

---

**Require:** Parameters  $\epsilon > 0$  and  $R > 0$ .

- 1: Initialize  $\mathbf{w}_1 = \mathbf{s}_1 = \mathbf{0}$ ,  $B_0 = \epsilon$ ,  $\mathbf{G}_0 = \mathbf{0}$ , and  $Q_0 = \epsilon^2$ .
  - 2: **for**  $t = 1, 2, \dots$  **do**
  - 3:   Play  $\mathbf{w}_t$  and observe  $\nabla_t \in \partial \ell_t(\mathbf{w}_t)$ .
  - 4:   Set  $B_t = B_{t-1} \vee \|\nabla_t\|$  and  $\bar{\nabla}_t := \nabla_t \cdot B_{t-1}/B_t$ .
  - 5:   Set  $\tilde{\nabla}_t = \bar{\nabla}_t - \mathbb{I}_{\langle \nabla_t, \mathbf{w}_t \rangle < 0} \cdot \langle \bar{\nabla}_t, \mathbf{s}_t \rangle \mathbf{s}_t$ . //  $\mathbf{s}_t$  is defined on Lines 1 and 8
  - 6:   Set  $\mathbf{G}_t = \mathbf{G}_{t-1} + \tilde{\nabla}_t$  and  $Q_t = Q_{t-1} + \|\tilde{\nabla}_t\|^2$ .
  - 7:   Set  $\mathbf{x}_{t+1} := -\mathbf{G}_t \cdot \frac{(2Q_t + B_t \|\mathbf{G}_t\|) \cdot \epsilon^2}{2(Q_t + B_t \|\mathbf{G}_t\|)^2 \sqrt{Q_t}} \cdot \exp\left(\frac{\|\mathbf{G}_t\|^2}{2Q_t + 2B_t \|\mathbf{G}_t\|}\right)$ .
  - 8:   Set  $\mathbf{w}_{t+1} = \Pi_{\mathcal{B}(R)}(\mathbf{x}_{t+1})$  and  $\mathbf{s}_{t+1} = \mathbf{x}_{t+1}/\|\mathbf{x}_{t+1}\| \cdot \mathbb{I}_{\|\mathbf{x}_{t+1}\| > R}$ .  
      // We use the convention that  $0/0 = 0$ .
  - 9: **end for**
- 

Another algorithm that will be instrumental to developing our projection-free (and scale-free) algorithms for strongly convex losses is FreeGrad (see Algorithm 9). An important property of



FreeGrad that will be useful to us when developing a projection-free algorithm for strongly convex losses (see Section C.2) is that its regret scales directly with the norm  $\|\mathbf{w}\|$  of the comparator, as opposed to the worst-case  $R$ . We note that FreeGrad internally clips the sequence of observed sub-gradients. In our application of FreeGrad, it will be useful to state the guarantee of its iterates  $(\mathbf{w}_t)$  in response to the sequence of clipped subgradients  $(\bar{\nabla}_t)$ :

**Proposition 26 (FreeGrad’s clipped linearized regret)** *For any adversarial sequence of convex losses  $(\ell_t)$  on  $\mathbb{R}^d$ , the iterates  $(\mathbf{w}_t)$  of FreeGrad with parameter  $\epsilon, R > 0$ , satisfy for all  $T \geq 1$  and  $\mathbf{w} \in \mathcal{B}(R)$ ,*

$$\sum_{t=1}^T \langle \bar{\nabla}_t, \mathbf{w}_t - \mathbf{w} \rangle \leq 2\|\mathbf{w}\| \sqrt{\bar{V}_T \ln_+ \left( \frac{2\|\mathbf{w}\|\bar{V}_T}{\epsilon^2} \right)} + 4B_T \|\mathbf{w}\| \ln \left( \frac{4B_T \|\mathbf{w}\| \sqrt{\bar{V}_T}}{\epsilon^2} \right) + \epsilon,$$

where  $\nabla_t \in \partial \ell_t(\mathbf{w}_t)$  (any sub-grad.),  $B_T := \epsilon \vee \max_{t \in [T]} \|\nabla_t\|$ ,  $\bar{\nabla}_t := \nabla_t \cdot B_{t-1}/B_t$ , and  $\bar{V}_T := \epsilon^2 + \sum_{t=1}^T \|\bar{\nabla}_t\|^2$ .

Technically, the original version of FreeGrad applies to unbounded OCO, whereas the version of FreeGrad in Algorithm 9 generates outputs in  $\mathcal{B}(R)$  using the same constrained-to-unconstrained reduction due to Cutkosky (2020). We constrain the outputs of FreeGrad to ensure that the iterates  $(\mathbf{w}_t)$  of Algorithm 1 in the setting of Section C.2 are bounded, which in turn ensures that the approximation errors of  $\text{OPT}_{\mathcal{C}^\circ}$  in Algorithm 1 are not too large. The proof of Proposition 26 follows by (Mhammedi and Koolen, 2020, Theorem 6) and (Cutkosky, 2020, Theorem 2).

We close this section by mentioning that the unbounded version of FreeGrad, where  $\tilde{\nabla}_t = \nabla_t$  and  $\mathbf{w}_t = \mathbf{x}_t$ , is an FTRL instance with the specific sequence of regularizers  $(\phi_t^*)$ , where  $\phi_t^*$  is the Fenchel dual of

$$\phi_t(\mathbf{x}) := \frac{1}{\sqrt{V_{t-1}}} \exp \left( \frac{\|\mathbf{x}\|^2}{2V_{t-1} + 2\|\mathbf{x}\|} \right), \quad V_t := \epsilon^2 + \sum_{s=1}^t \|\nabla_s\|^2.$$

In particular, the iterate  $\mathbf{x}_{t+1}$  on Line 7 of Algorithm 9 is given by

$$\mathbf{x}_{t+1} \in \arg \min_{\mathbf{x} \in \mathbb{R}^d} \{ \langle \mathbf{G}_t, \mathbf{x} \rangle + \phi_t^*(\mathbf{x}) \} = \{ \nabla \phi_t(-\mathbf{G}_t) \}.$$

We recall that the Fenchel dual  $\phi_t^*$  of  $\phi_t$  is defined as  $\phi_t^*(\mathbf{x}) = \sup_{\mathbf{u} \in \mathbb{R}^d} \{ \langle \mathbf{u}, \mathbf{x} \rangle - \phi_t(\mathbf{u}) \}$ .

## Appendix H. Linear Optimization on $\mathcal{C}^\circ$ using a Membership Oracle for $\mathcal{C}$

In this section, we restate and prove a slight extension of (Lee et al., 2018, Lemma 9 & 10), which we need in the proof of Proposition 11 at the end of this section. Our extension involves showing that the approximate subgradient in (Lee et al., 2018, Lemma 10) has bounded norm in high probability, which we need in the proof of Proposition 11. We also make the limiting argument used in the proof of (Lee et al., 2018, Lemma 9 & 10) more explicit; in particular, for any convex function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  (not necessarily differentiable) and  $\epsilon > 0$ , we explicitly construct a twice differentiable function  $\tilde{f}_\epsilon$  such that  $\|f - \tilde{f}_\epsilon\|_\infty \leq \epsilon$ . This will then allow us to make a limiting argument precise via Fatou’s lemma to get the final result we want (see proof of Lemma 29).

Throughout this section, we let  $U_\infty(\mathbf{u}, \nu)$  denote the uniform distribution over  $\mathcal{B}_\infty(\mathbf{u}, \nu)$ , for any  $\mathbf{u} \in \mathbb{R}^d$  and  $\nu > 0$ . The next lemma is taken from (Lee et al., 2018, Lemma 9) with only minor notation adjustments.



**Lemma 27** For any  $\mathbf{w} \in \mathbb{R}^d$ ,  $0 < \nu_2 \leq \nu_1$ , and twice differentiable convex function  $h$  defined on  $B_\infty(\mathbf{w}, \nu_1 + \nu_2)$  with  $\|\nabla h(\mathbf{z})\|_\infty \leq L$  for any  $\mathbf{z} \in B_\infty(\mathbf{w}, \nu_1 + \nu_2)$  we have

$$\mathbb{E}_{\mathbf{u} \sim U_\infty(\mathbf{w}, \nu_1)} \mathbb{E}_{\mathbf{z} \sim U_\infty(\mathbf{u}, \nu_2)} \|\nabla h(\mathbf{z}) - \mathbf{g}(\mathbf{u})\|_1 \leq \frac{\nu_2 d^{3/2} L}{\nu_1},$$

where  $\mathbf{g}(\mathbf{u}) := \mathbb{E}_{\mathbf{z} \sim U_\infty(\mathbf{u}, \nu_2)}[\nabla h(\mathbf{z})]$  and  $U_\infty(\mathbf{u}, \nu)$  denotes the uniform distribution over  $B_\infty(\mathbf{u}, \nu)$ .

---

**Algorithm 10** Approximate Subgradient Oracle (Lee et al., 2018).

---

**Require:** Inputs  $\nu_1 > 0$ ,  $L > 0$ ,  $\varepsilon > 0$ , and  $\mathbf{w} \in \mathbb{R}^d$ .

- 1: Function  $\tilde{f} : \mathbb{R}^d \rightarrow \mathbb{R}$ .
  - 2: Set  $\nu_2 = \sqrt{\frac{\varepsilon \nu_1}{d^{1/2} L}}$ .
  - 3: Sample  $\mathbf{u} \in B_\infty(\mathbf{w}, \nu_1)$  and  $\mathbf{z} \in B_\infty(\mathbf{u}, \nu_2)$  independently and uniformly at random.
  - 4: **for**  $i = 1, 2, \dots, d$  **do**
  - 5:   Let  $\mathbf{w}'_i$  and  $\mathbf{w}_i$  be the end point of the interval  $B_\infty(\mathbf{u}, \nu_2) \cap \{\mathbf{z} + \lambda \mathbf{e}_i : \lambda \in \mathbb{R}\}$ .
  - 6:   Set  $\tilde{s}_i = \frac{1}{2\nu_2}(\tilde{f}(\mathbf{w}'_i) - \tilde{f}(\mathbf{w}_i))$ .
  - 7: **end for**
  - 8: Set  $\tilde{\mathbf{s}} = (\tilde{s}_1, \dots, \tilde{s}_d)^\top$ .
  - 9: Return  $\tilde{\mathbf{s}}$ .
- 

We now restate and slightly extend (Lee et al., 2018, Lemma 10) for twice differentiable functions. We then extend the result to non-differentiable functions using Fatou’s lemma—see Lemma 29.

**Lemma 28** Let  $L, \nu_1 > 0$ ,  $\mathbf{w} \in \mathbb{R}^d$ , and  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  be a twice differentiable convex function such that  $\|\nabla h(\mathbf{x})\|_\infty \leq L$ , for any  $\mathbf{x} \in B_\infty(\mathbf{w}, 2\nu_1)$ . Also, let  $\varepsilon \in (0, \nu_1 \sqrt{d}L]$  and  $\tilde{h} : \mathbb{R}^d \rightarrow \mathbb{R}$  be such that  $\|\tilde{h} - h\|_\infty \leq \varepsilon'$  for some  $\varepsilon' > 0$ . Then, the variables  $\nu_2$ ,  $\mathbf{u}$ ,  $\mathbf{z}$ ,  $\tilde{\mathbf{s}}$ , and  $(\mathbf{w}_i, \mathbf{w}'_i)_{i \in [d]}$  generated during a run of Alg. 10 with input  $(\tilde{h}, \nu_1, L, \varepsilon, \mathbf{w})$  satisfy

$$\forall \mathbf{v} \in \mathbb{R}^d, h(\mathbf{v}) \geq h(\mathbf{w}) + \langle \tilde{\mathbf{s}}, \mathbf{v} - \mathbf{w} \rangle - \|\nabla h(\mathbf{z}) - \tilde{\mathbf{s}}\|_1 \cdot \|\mathbf{v} - \mathbf{w}\|_\infty - 4d\nu_1 L,$$

Furthermore, for  $\tilde{s}_i := (h(\mathbf{w}'_i) - h(\mathbf{w}_i))/(2\nu_2)$ , for all  $i \in [d]$ , we have

$$\|\nabla h(\mathbf{z}) - \tilde{\mathbf{s}}\|_1 \leq \|\nabla h(\mathbf{z}) - \tilde{\mathbf{s}}\|_1 + d\varepsilon'/\nu_2 \quad \text{and} \quad \mathbb{E}[\|\nabla h(\mathbf{z}) - \tilde{\mathbf{s}}\|_1] \leq 2d^{5/4} \sqrt{L\varepsilon/\nu_1}.$$

**Proof** Let  $\mathbf{u}$ ,  $\mathbf{w}_i$ , and  $\mathbf{w}'_i$  be the random vectors generated during the call to Algorithm 10 in the lemma’s statement. Further, let  $\tilde{\mathbf{s}} \in \mathbb{R}^d$  be such that  $\tilde{s}_i := (h(\mathbf{w}'_i) - h(\mathbf{w}_i))/(2\nu_2)$ , for  $i \in [d]$ . Applying the convexity of  $h$  yields, for any  $\mathbf{v} \in \mathbb{R}^d$ ,

$$\begin{aligned} h(\mathbf{v}) &\geq h(\mathbf{z}) + \langle \nabla h(\mathbf{z}), \mathbf{v} - \mathbf{z} \rangle \\ &= h(\mathbf{z}) + \langle \tilde{\mathbf{s}}, \mathbf{v} - \mathbf{w} \rangle + \langle \nabla h(\mathbf{z}) - \tilde{\mathbf{s}}, \mathbf{v} - \mathbf{w} \rangle + \langle \nabla h(\mathbf{z}), \mathbf{w} - \mathbf{z} \rangle \\ &\geq h(\mathbf{w}) + \langle \tilde{\mathbf{s}}, \mathbf{v} - \mathbf{w} \rangle - \|\nabla h(\mathbf{z}) - \tilde{\mathbf{s}}\|_1 \|\mathbf{v} - \mathbf{w}\|_\infty - \|\nabla h(\mathbf{z})\|_\infty \|\mathbf{w} - \mathbf{z}\|_1. \end{aligned}$$

Now,  $\|\nabla h(\mathbf{z})\|_\infty \leq L$  and  $\|\mathbf{w} - \mathbf{z}\|_1 \leq d \cdot \|\mathbf{w} - \mathbf{z}\|_\infty \leq 2d(\nu_1 + \nu_2)$  by definition of  $\mathbf{z}$  in Algorithm 10. Furthermore, since  $\varepsilon \leq \nu_1 \sqrt{d}L$ , we have  $\nu_2 = \sqrt{\frac{\varepsilon \nu_1}{d^{1/2} L}} \leq \nu_1$ . Plugging these facts in the inequality of the previous display implies

$$h(\mathbf{v}) \geq h(\mathbf{w}) + \langle \tilde{\mathbf{s}}, \mathbf{v} - \mathbf{w} \rangle - \|\nabla h(\mathbf{z}) - \tilde{\mathbf{s}}\|_1 \|\mathbf{v} - \mathbf{w}\|_\infty - 4d\nu_1 L.$$

This shows the first inequality we are after. Now, by the definition of  $\tilde{\mathbf{s}}$  in Algorithm 10 and the fact that  $\tilde{h}$  satisfies  $\|\tilde{h} - h\|_\infty \leq \varepsilon'$ , we have

$$\|\nabla h(\mathbf{z}) - \tilde{\mathbf{s}}\|_1 \leq \|\nabla h(\mathbf{z}) - \tilde{\mathbf{s}}\|_1 + d\varepsilon'/\nu_2.$$

It remains to bound  $\mathbb{E}[\|\nabla h(\mathbf{z}) - \tilde{\mathbf{s}}\|]$ . For this, let  $\mathbf{g}(\mathbf{u}) := \mathbb{E}_{\mathbf{z} \sim U_\infty(\mathbf{u}, \nu_2)}[\nabla h(\mathbf{z})]$  and note that

$$\begin{aligned} \mathbb{E}_{\mathbf{z}} [|\tilde{s}_i - [\mathbf{g}(\mathbf{u})]_i|] &= \mathbb{E}_{\mathbf{z}} \left[ \left| \frac{h(\mathbf{w}'_i) - h(\mathbf{w}_i)}{2\nu_2} - [\mathbf{g}(\mathbf{u})]_i \right| \right], \\ &\leq \mathbb{E}_{\mathbf{z}} \left[ \frac{1}{2\nu_2} \int \left| \frac{dh}{dw_i}(\mathbf{z} + \lambda \mathbf{e}_i) - [\mathbf{g}(\mathbf{u})]_i \right| d\lambda \right], \\ &= \mathbb{E}_{\mathbf{z}} \left[ \left| \frac{dh}{dw_i}(\mathbf{z}) - [\mathbf{g}(\mathbf{u})]_i \right| \right], \end{aligned}$$

where we used that both  $\mathbf{z} + \lambda \mathbf{e}_i$  and  $\mathbf{z}$  are uniform distribution on  $\mathcal{B}_\infty(\mathbf{u}, \nu_2)$  in the last line. Hence, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{z}} [\|\tilde{\mathbf{s}} - \nabla h(\mathbf{z})\|_1] &\leq \mathbb{E}_{\mathbf{z}} [\|\nabla h(\mathbf{z}) - \mathbf{g}(\mathbf{u})\|_1] + \mathbb{E}_{\mathbf{z}} \|\tilde{\mathbf{s}} - \mathbf{g}(\mathbf{u})\|_1 \leq 2\mathbb{E}_{\mathbf{z}} [\|\nabla h(\mathbf{z}) - \mathbf{g}(\mathbf{u})\|_1] \\ &\leq 2d^{5/4} \sqrt{\frac{\varepsilon L}{\nu_1}}, \end{aligned}$$

where the last inequality follows by Lemma 27 and the fact that  $\nu_2 = \sqrt{\frac{\varepsilon \nu_1}{d^{1/2} L}} \leq \nu_1$ .  $\blacksquare$

**Lemma 29** *Let  $L, \nu_1 > 0$ ,  $\mathbf{w} \in \mathbb{R}^d$ , and  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex (not necessarily differentiable) function such that  $\sup_{\mathbf{g} \in \partial f(\mathbf{x})} \|\mathbf{g}\| \leq L$ , for any  $\mathbf{x} \in B_\infty(\mathbf{w}, 2\nu_1)$ . Also, let  $\varepsilon \in (0, \nu_1 \sqrt{d} L)$  and  $\tilde{f} : \mathbb{R}^d \rightarrow \mathbb{R}$  be such that  $\|\tilde{f} - f\|_\infty \leq \varepsilon$ . Then, the output  $\tilde{\mathbf{s}}$  of Algorithm 10 with input  $(\tilde{f}, \nu_1, L, \varepsilon, \mathbf{w})$  satisfies*

$$\begin{aligned} \forall \mathbf{v} \in \mathbb{R}^d, f(\mathbf{v}) &\geq f(\mathbf{w}) + \langle \tilde{\mathbf{s}}, \mathbf{v} - \mathbf{w} \rangle - X \cdot \|\mathbf{v} - \mathbf{w}\|_\infty - 4d\nu_1 L; \\ \|\tilde{\mathbf{s}}\|_\infty &\leq L + d^{1/4} \sqrt{L\varepsilon/\nu_1}; \quad \|\tilde{\mathbf{s}}\| \leq X + L \leq \|\tilde{\mathbf{s}}\|_1 + (d+1)L; \\ \text{and } \|\tilde{\mathbf{s}}\|^2 &\leq \left(4L + d^{1/4} \sqrt{L\varepsilon/\nu_1}\right) X + L^2, \end{aligned}$$

where  $X \geq 0$  is a non-negative random variable satisfying  $\mathbb{E}[X] \leq 3d^{5/4} \sqrt{L\varepsilon/\nu_1}$ .

**Proof** Let  $\nu_2, \mathbf{u}, \mathbf{z}, \tilde{\mathbf{s}}$ , and  $(\mathbf{w}_i, \mathbf{w}'_i)_{i \in [d]}$  be the variables generated during the call to Algorithm 10 in the lemma's statement. Further, let  $\sigma > 0$  and  $h_\sigma : \mathbb{R}^d \rightarrow \mathbb{R}$  be such that

$$h_\sigma(\mathbf{w}) := \mathbb{E}[f(\mathbf{w} + \mathbf{x})], \quad \text{where } \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_d).$$

It is known that  $h_\sigma$  is twice differentiable (Bhatnagar, 2007; Nesterov and Spokoiny, 2017; Abernethy et al., 2014), and satisfies (Duchi et al., 2012, Lemma 9),

$$\forall \mathbf{w} \in \mathbb{R}^d, \quad \|\nabla h_\sigma(\mathbf{w})\| \leq L, \quad \text{and} \quad f(\mathbf{w}) \leq h_\sigma(\mathbf{w}) \leq f(\mathbf{w}) + L\sigma\sqrt{d}. \quad (49)$$

Thus, we have  $\|\tilde{f} - h_\sigma\|_\infty \leq \varepsilon_\sigma := \varepsilon + L\sigma\sqrt{d}$ , and so by Lemma 28

$$h_\sigma(\mathbf{v}) \geq h_\sigma(\mathbf{w}) + \langle \tilde{\mathbf{s}}, \mathbf{v} - \mathbf{w} \rangle - \|\nabla h_\sigma(\mathbf{z}) - \tilde{\mathbf{s}}\|_1 \cdot \|\mathbf{v} - \mathbf{w}\|_\infty - 4d\nu_1 L, \quad \forall \mathbf{v} \in \mathbb{R}^d,$$

and for  $\tilde{\mathbf{s}}_\sigma := (h_\sigma(\mathbf{w}'_1) - h_\sigma(\mathbf{w}_1), \dots, h_\sigma(\mathbf{w}'_d) - h_\sigma(\mathbf{w}_d))^\top / (2\nu_2)$ ,

$$\|\nabla h_\sigma(\mathbf{z}) - \tilde{\mathbf{s}}\|_1 \leq \|\nabla h_\sigma(\mathbf{z}) - \tilde{\mathbf{s}}_\sigma\|_1 + d\varepsilon_\sigma/\nu_2, \quad \text{and} \quad \mathbb{E}[\|\nabla h_\sigma(\mathbf{z}) - \tilde{\mathbf{s}}_\sigma\|_1] \leq 2d^{5/4}\sqrt{L\varepsilon/\nu_1}, \quad (50)$$

where the expectation is with respect to  $\mathbf{z} \sim U_\infty(\mathbf{u}, \nu_2)$  and  $\mathbf{u} \sim U_\infty(\mathbf{w}, \nu_1)$ . Combining this with (49) leads to

$$f(\mathbf{v}) \geq f(\mathbf{w}) + \langle \tilde{\mathbf{s}}, \mathbf{v} - \mathbf{w} \rangle - \|\nabla h_\sigma(\mathbf{z}) - \tilde{\mathbf{s}}\|_1 \cdot \|\mathbf{v} - \mathbf{w}\|_\infty - 4d\nu_1 L - L\sigma\sqrt{d}, \quad \forall \mathbf{v} \in \mathbb{R}^d. \quad (51)$$

Now, define  $X_n := \|\nabla h_{1/n}(\mathbf{z}) - \tilde{\mathbf{s}}\|_1$ ,  $n \geq 1$  and note that

$$0 \leq X_n \leq dL + \|\tilde{\mathbf{s}}\|_1 < +\infty, \quad \forall n.$$

Therefore, we have  $X := \liminf_{n \rightarrow \infty} X_n \leq \|\tilde{\mathbf{s}}\|_1 + dL < +\infty$ . Furthermore, by Fatou's lemma and (50), we have

$$\mathbb{E}[X] = \mathbb{E}[\liminf_{n \rightarrow \infty} X_n] \leq \liminf_{n \rightarrow \infty} \mathbb{E}[X_n] \leq \liminf_{n \rightarrow \infty} \mathbb{E}[\|\nabla h_{1/n}(\mathbf{z}) - \tilde{\mathbf{s}}_{1/n}\|_1] + d\varepsilon/\nu_2 \leq 3d^{5/4}\sqrt{L\varepsilon/\nu_1},$$

where the last inequality also follows by (50) and the fact that  $\nu_2 = \sqrt{\frac{\varepsilon\nu_1}{d^{1/2}L}}$ . Combining this with (51) leads to the first inequality of the lemma. Now, we have, for any  $i \in [d]$ ,

$$\tilde{s}_i = \frac{\tilde{f}(\mathbf{w}'_i) - \tilde{f}(\mathbf{w}_i)}{2\nu_2} \leq \frac{|f(\mathbf{w}'_i) - f(\mathbf{w}_i)|}{2\nu_2} + \frac{\varepsilon}{\nu_2} \leq L + \frac{\varepsilon}{\nu_2} = L + d^{1/4}\sqrt{\frac{L\varepsilon}{\nu_1}},$$

where the last inequality follows by the definition of  $\mathbf{w}'_i$  and  $\mathbf{w}_i$ , and the fact that  $\sup_{\mathbf{g} \in \partial f(\mathbf{x})} \|\mathbf{g}\| \leq L$  (and so  $f$  is  $L$ -Lipschitz). In particular, this implies that

$$\forall i \in [d], \quad |\tilde{s}_i - [\nabla h_{1/n}(\mathbf{z})]_i| \leq 2L + d^{1/4}\sqrt{L\varepsilon/\nu_1}. \quad (52)$$

Using this, we get

$$\begin{aligned} \|\tilde{\mathbf{s}}\|^2 &\leq \|\tilde{\mathbf{s}} - \nabla h_{1/n}(\mathbf{z})\|^2 + 2\|\tilde{\mathbf{s}} - \nabla h_{1/n}(\mathbf{z})\| \|\nabla h_{1/n}(\mathbf{z})\| + \|\nabla h_{1/n}(\mathbf{z})\|^2, \\ &\leq \left(2L + d^{1/4}\sqrt{L\varepsilon/\nu_1}\right) \|\tilde{\mathbf{s}} - \nabla h_{1/n}(\mathbf{z})\|_1 + 2\|\tilde{\mathbf{s}} - \nabla h_{1/n}(\mathbf{z})\|_1 L + L^2, \quad (\text{by (52)}) \\ &\leq \left(4L + d^{1/4}\sqrt{L\varepsilon/\nu_1}\right) \|\tilde{\mathbf{s}} - \nabla h_{1/n}(\mathbf{z})\|_1 + L^2, \\ &\leq \left(4L + d^{1/4}\sqrt{L\varepsilon/\nu_1}\right) \liminf_{n \rightarrow \infty} \|\tilde{\mathbf{s}} - \nabla h_{1/n}(\mathbf{z})\|_1 + L^2, \\ &= \left(4L + d^{1/4}\sqrt{L\varepsilon/\nu_1}\right) X + L^2. \end{aligned}$$

It remains to bound the norm  $\|\tilde{\mathbf{s}}\|$ . Similar to how we bounded  $\|\tilde{\mathbf{s}}\|^2$ , we have

$$\begin{aligned} \|\tilde{\mathbf{s}}\| &\leq \|\tilde{\mathbf{s}} - \nabla h_{1/n}(\mathbf{z})\| + \|\nabla h_{1/n}(\mathbf{z})\|, \\ &\leq \|\tilde{\mathbf{s}} - \nabla h_{1/n}(\mathbf{z})\|_1 + L, \\ &\leq \liminf_{n \rightarrow \infty} \|\tilde{\mathbf{s}} - \nabla h_{1/n}(\mathbf{z})\|_1 + L = X + L. \end{aligned}$$

This completes the proof. ■

### H.1. Proof of Lemma 10 (Approximate Gauge Function)

**Proof of Lemma 10.** Let  $\epsilon = \delta r / (4\kappa)^2$ . We will first show that for all  $\mathbf{w} \in \mathcal{B}(6R/5)$ ,

$$[\gamma_{\mathcal{C}}(\mathbf{w}) \geq 9/16 \text{ or } \tilde{\gamma} \geq 1] \implies [\text{MEM}_{\mathcal{C}}(2\mathbf{w}; \epsilon) = 0 \text{ and } \|\mathbf{w}\| \geq r/2]. \quad (53)$$

Suppose that  $\tilde{\gamma} \geq 1$ . By Lines 2 and 3 of Algorithm 2, this implies that  $\text{MEM}_{\mathcal{C}}(2\mathbf{w}; \epsilon) = 0$ . Now suppose that  $\gamma_{\mathcal{C}}(\mathbf{w}) \geq 9/16$ . Note that  $\text{MEM}_{\mathcal{C}}(2\mathbf{w}; \epsilon) = 1$  only if  $2\mathbf{w} \in \mathcal{B}(\mathcal{C}, \epsilon)$ . By the fact that  $\mathcal{B}(r) \subseteq \mathcal{C}$  (Assumption 1) and  $\delta < 1$ , we have that  $\mathcal{C} + \mathcal{B}(\epsilon) = \mathcal{B}(\mathcal{C}, \epsilon)$ . Thus, by a standard result in convex analysis, see e.g. (Hiriart-Urruty and Lemaréchal, 2004, Thm C.3.3.2), we have, for all  $\mathbf{x} \in \mathcal{C}$ ,

$$\sigma_{\mathcal{B}(\mathcal{C}, \epsilon)}(\mathbf{x}) = \sigma_{\mathcal{C}}(\mathbf{x}) + \sigma_{\mathcal{B}(\epsilon)}(\mathbf{x}). \quad (54)$$

Let  $\mathbf{x}_* \in \partial\gamma_{\mathcal{C}}(\mathbf{w}) = \arg \max_{\mathbf{x} \in \mathcal{C}^\circ} \langle \mathbf{x}, \mathbf{w} \rangle$ . Since  $\mathbf{x}_* \in \mathcal{C}^\circ$ , we have  $\langle \mathbf{x}_*, \mathbf{y} \rangle \leq 1$ , for all  $\mathbf{y} \in \mathcal{C}$ , by definition of  $\mathcal{C}^\circ$ . Further, since  $\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) \in \mathcal{C}$  (implied by Lemma 6) and  $\langle \mathbf{x}_*, \mathbf{w} \rangle = \gamma_{\mathcal{C}}(\mathbf{w})$ , we have

$$\sigma_{\mathcal{C}}(\mathbf{x}_*) = \sup_{\mathbf{y} \in \mathcal{C}} \langle \mathbf{x}_*, \mathbf{y} \rangle = \langle \mathbf{x}_*, \mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) \rangle = 1. \quad (55)$$

Plugging this into (54) and using the fact that  $\sigma_{\mathcal{B}(\epsilon)}(\mathbf{x}_*) = \delta r \|\mathbf{x}_*\| / (8\kappa^2)$ , we get

$$\sigma_{\mathcal{B}(\mathcal{C}, \epsilon)}(\mathbf{x}_*) = 1 + \delta r \|\mathbf{x}_*\| / (8\kappa^2) \leq 1 + \delta / (8\kappa^2) < 9/8, \quad (56)$$

where the last inequality follows by the fact that  $\delta \in (0, 1)$  and  $\|\mathbf{x}_*\| \leq 1/r$  (because  $\mathbf{x}_* \in \mathcal{C}^\circ$ , see Lemma 42-(c)). On the other hand, since  $\gamma_{\mathcal{C}}(\mathbf{w}) \geq 9/16$ , we have, by (55)

$$9/8 = 9/8 \langle \mathbf{x}_*, \mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) \rangle \leq 9/8 \langle \mathbf{x}_*, 16\mathbf{w}/9 \rangle = \langle \mathbf{x}_*, 2\mathbf{w} \rangle.$$

This inequality together with (56) implies that the vector  $\mathbf{x}_*$  separates  $2\mathbf{w}$  from  $\mathcal{B}(\mathcal{C}, \epsilon)$  and so we have  $\text{MEM}_{\mathcal{C}}(2\mathbf{w}; \epsilon) = 0$  by definition of the Membership Oracle  $\text{MEM}_{\mathcal{C}}(\cdot; \epsilon)$ . It remains to show that  $\|\mathbf{w}\| \geq r/2$ . If  $\gamma_{\mathcal{C}}(\mathbf{w}) \geq 9/16$ , then  $\|\mathbf{w}\| \geq 9r/16$  by the fact that  $\gamma_{\mathcal{C}}(\mathbf{w}) \leq \|\mathbf{w}\|/r$  (Lemma 42-(c)). Also, if  $\tilde{\gamma} \geq 1$ , then Lines 2 and 3 of Algorithm 2 imply that  $\|\mathbf{w}\| \geq r/2$ . So far, we have shown that (53) holds. By definition of  $\alpha$  and  $\beta$  in Algorithm 2, (53) implies that if  $\gamma_{\mathcal{C}}(\mathbf{w}) \geq 9/16$  or  $\tilde{\gamma} \geq 1$ , then

$$\text{MEM}_{\mathcal{C}}(\alpha\mathbf{w}; \epsilon) = 1 \text{ and } \text{MEM}_{\mathcal{C}}(\beta\mathbf{w}; \epsilon) = 0. \quad (57)$$

Next, we will show that if either  $\gamma_{\mathcal{C}}(\mathbf{w}) \geq 9/16$  or  $\tilde{\gamma} \geq 1$ , then for any  $\mu > 0$ ,

$$\text{MEM}_{\mathcal{C}}(\mu\mathbf{w}; \epsilon) = 0 \implies \mu \geq \frac{1}{\gamma_{\mathcal{C}}(\mathbf{w})} - \frac{\delta}{8\kappa^2} \text{ and } \text{MEM}_{\mathcal{C}}(\mu\mathbf{w}; \epsilon) = 1 \implies \mu \leq \frac{1}{\gamma_{\mathcal{C}}(\mathbf{w})} + \frac{\delta}{8\kappa^2}. \quad (58)$$

Suppose that  $\text{MEM}_{\mathcal{C}}(\mu\mathbf{w}; \epsilon) = 0$ . By (2), this implies that  $\mu\mathbf{w} \notin \mathcal{B}(\mathcal{C}, -\epsilon)$ . On the other hand, since  $\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w})$  is on the boundary of  $\mathcal{C}$  (by definition of the Gauge function and Lemma 6), the vector  $\mathbf{v} := \mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) - \epsilon\mathbf{w}/\|\mathbf{w}\|$  is in  $\mathcal{B}(\mathcal{C}, -\epsilon)$ . Since  $\mathbf{v}$  and  $\mu\mathbf{w}$  are aligned, the fact that  $\mu\mathbf{w} \notin \mathcal{B}(\mathcal{C}, -\epsilon)$  and  $\mathcal{C}$  is convex implies that

$$\|\mu\mathbf{w}\| \geq \|\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) - \epsilon\mathbf{w}/\|\mathbf{w}\|\| = \|\mathbf{w}\|/\gamma_{\mathcal{C}}(\mathbf{w}) - r\delta/(16\kappa^2) \geq \|\mathbf{w}\|/\gamma_{\mathcal{C}}(\mathbf{w}) - \|\mathbf{w}\|\delta/(8\kappa^2),$$

where the last inequality follows by the fact that  $\|\mathbf{w}\| \geq r/2$  (see (53)). Dividing by  $\|\mathbf{w}\|$  (this is non-zero by (53)) on both sides shows the first implication in (58). Similarly, if  $\text{MEM}_{\mathcal{C}}(\mu\mathbf{w}; \epsilon) = 1$ . Again by (2), this implies that  $\mu\mathbf{w} \in \mathcal{B}(\mathcal{C}, \epsilon)$ . Combining this with the fact that  $\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) + \epsilon\mathbf{w}/\|\mathbf{w}\|$  is on the boundary of  $\mathcal{B}(\mathcal{C}, \epsilon)$  (since  $\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w})$  is on the boundary of  $\mathcal{C}$ ) implies that

$$\|\mu\mathbf{w}\| \leq \|\mathbf{w}/\gamma_{\mathcal{C}}(\mathbf{w}) + \epsilon\mathbf{w}/\|\mathbf{w}\|\| = \|\mathbf{w}\|/\gamma_{\mathcal{C}}(\mathbf{w}) + \delta r/(16\kappa^2) \leq \|\mathbf{w}\|/\gamma_{\mathcal{C}}(\mathbf{w}) + \delta\|\mathbf{w}\|/(8\kappa^2),$$

where the last inequality follows by the fact that  $\|\mathbf{w}\| \geq r/2$  (see (53)). This shows the second implication in (58). Let  $\tilde{\alpha}$  and  $\tilde{\beta}$  be the values of  $\alpha$  and  $\beta$  in Algorithm 2 after the while-loop is over. By combining (57), which we recall holds when  $\gamma_{\mathcal{C}}(\mathbf{w}) \geq 9/16$  or  $\tilde{\gamma} \geq 1$ , and (58), we get

$$\tilde{\gamma} \geq \gamma_{\mathcal{C}}(\mathbf{w}); \text{ and } \frac{1}{\gamma_{\mathcal{C}}(\mathbf{w})} - \frac{1}{\tilde{\gamma}} = \frac{1}{\gamma_{\mathcal{C}}(\mathbf{w})} - \left( \tilde{\alpha} - \frac{\delta}{8\kappa^2} \right) \leq \tilde{\beta} - \tilde{\alpha} + \frac{\delta}{8\kappa} + \frac{\delta}{8\kappa^2} \leq \frac{3\delta}{8\kappa^2}, \quad (59)$$

where the last inequality follows by the fact that the while-loop in Algorithm 2 terminates when  $\beta - \alpha \leq \delta/(8\kappa^2)$ . Multiplying both sides of (59) by  $\tilde{\gamma} \cdot \gamma_{\mathcal{C}}(\mathbf{w})$ , we get

$$\tilde{\gamma} - \gamma_{\mathcal{C}}(\mathbf{w}) \leq \tilde{\gamma} \cdot \gamma_{\mathcal{C}}(\mathbf{w}) \cdot 3\delta/(8\kappa^2). \quad (60)$$

Using (59), we get that  $\tilde{\gamma} \leq (1/\gamma_{\mathcal{C}}(\mathbf{w}) - 3\delta/(8\kappa^2))^{-1}$ . Plugging this into (60), we get

$$\tilde{\gamma} - \gamma_{\mathcal{C}}(\mathbf{w}) \leq \frac{\gamma_{\mathcal{C}}(\mathbf{w}) \cdot 3\delta/(8\kappa^2)}{1/\gamma_{\mathcal{C}}(\mathbf{w}) - 3\delta/(8\kappa^2)} \stackrel{(a)}{\leq} \frac{6\kappa/5 \cdot 3\delta/(8\kappa^2)}{5/(6\kappa) - 3\delta/(8\kappa^2)} \leq \delta,$$

where (a) follows by the fact that  $\gamma_{\mathcal{C}}(\mathbf{w}) \leq \|\mathbf{w}\|/r \leq 6\kappa/5$ , for all  $\mathbf{w} \in \mathcal{B}(6R/5)$ , and the last inequality follows by the fact that  $\delta \in (0, 1)$  and  $\kappa \geq 1$ .

We now consider the computational complexity. Note that at every iteration of the while-loop in Line 7 of Algorithm 2, the difference  $\beta - \alpha$  halves. Thus, since  $\beta - \alpha$  is initially equal to 2, the algorithm terminates (i.e. when  $\beta - \alpha \leq \delta/(8\kappa^2)$ ) after at most  $\lceil \log_2((4\kappa)^2/\delta) \rceil + 1$  steps.  $\blacksquare$

## H.2. Proof of Proposition 11 (Approximate LO Oracle on $\mathcal{C}^\circ$ )

**Proof of Proposition 11.** Let  $\epsilon$  and  $\tilde{\gamma}$  be as in Algorithm 3. First, suppose that  $\tilde{\gamma} \geq 1$ . We will show (17) using the result of Lemma 29 with  $(f, \tilde{f}) = (\gamma_{\mathcal{C}}, \text{GAU}_{\mathcal{C}}(\cdot; \epsilon))$  and  $(\nu_1, \nu_2, \mathbf{w})$  as in Algorithm 3. First, note that by our choice of  $\nu_1$  and  $\epsilon$ , the technical condition  $\epsilon \leq \nu_1 \sqrt{d}L$  under which Lemma 29 holds translates to  $\delta \leq 10d^{3/2}\kappa$ . This holds since  $\delta \in (0, 1)$  and  $\kappa, d \geq 1$ . We also need to check the technical condition  $\|f - \tilde{f}\|_{\infty} \leq \epsilon$  of Lemma 29. We note that in Algorithm 3 we only ever evaluate  $\tilde{f} = \text{GAU}_{\mathcal{C}}(\cdot; \epsilon)$  at the points  $(\mathbf{w}_i, \mathbf{w}'_i)_{i \in [d]}$ , and so we only need to check the condition

$$|\gamma_{\mathcal{C}}(\mathbf{w}_i) - \text{GAU}_{\mathcal{C}}(\mathbf{w}_i; \epsilon)| \vee |\gamma_{\mathcal{C}}(\mathbf{w}'_i) - \text{GAU}_{\mathcal{C}}(\mathbf{w}'_i; \epsilon)| \leq \epsilon, \quad \forall i \in [d]. \quad (61)$$

The condition would follow from Lemma 10 if  $\mathbf{w}_i$  and  $\mathbf{w}'_i$  satisfy  $\gamma_{\mathcal{C}}(\mathbf{w}_i) \wedge \gamma_{\mathcal{C}}(\mathbf{w}'_i) \geq 9/16$  and  $\mathbf{w}_i, \mathbf{w}'_i \in \mathcal{B}(6R/5)$ , for all  $i \in [d]$ . Let  $i \in [d]$  and  $\mathbf{v} \in \{\mathbf{w}_i, \mathbf{w}'_i\}$ . By definition of  $\mathbf{w}_i, \mathbf{w}'_i$ , and  $\mathbf{u}$  in Algorithm 3, we have

$$\|\mathbf{w} - \mathbf{v}\| \leq \sqrt{d}\|\mathbf{w} - \mathbf{u}\|_{\infty} + \sqrt{d}\|\mathbf{u} - \mathbf{v}\|_{\infty} \leq \sqrt{d}\nu_1 + \sqrt{d}\nu_2 \leq \frac{11r\delta}{100}, \quad (62)$$

where the last inequality follows by our choice of  $\nu_1$  and  $\nu_2$ . Combining (62) with the fact that  $\mathbf{w} \in \mathcal{B}(R)$  and triangular inequality, we get  $\|\mathbf{v}\| \leq \|\mathbf{w}\| + 11r\delta/100 \leq 6R/5$ , and so  $\mathbf{v} \in \mathcal{B}(6R/5)$ . On the other hand, by the fact that  $\gamma_{\mathcal{C}} = \sigma_{\mathcal{C}^\circ}$  and the sub-additivity of the support function (Lemma 42-(e)), we have

$$\gamma_{\mathcal{C}}(\mathbf{v}) = \sigma_{\mathcal{C}^\circ}(\mathbf{v}) \geq \sigma_{\mathcal{C}^\circ}(\mathbf{w}) - \sigma_{\mathcal{C}^\circ}(\mathbf{w} - \mathbf{v}) = \gamma_{\mathcal{C}}(\mathbf{w}) - \gamma_{\mathcal{C}}(\mathbf{w} - \mathbf{v}) \stackrel{(a)}{\geq} \gamma_{\mathcal{C}}(\mathbf{w}) - \frac{11\delta}{100} \geq \frac{9}{16},$$

where (a) follows by the fact that  $\gamma_{\mathcal{C}}(\mathbf{w} - \mathbf{v}) \leq \|\mathbf{w} - \mathbf{v}\|/r$  (Lemma 42-(c)) and (62), and the last inequality follows by the fact that  $\delta \leq 1/3$  and  $\gamma_{\mathcal{C}}(\mathbf{w}) \geq 2/3$ ; the latter is because  $\tilde{\gamma} \geq 1$  and  $\gamma_{\mathcal{C}}(\mathbf{w}) \geq \tilde{\gamma} - \delta$  (Lemma 10). Therefore, by Lemma 10, (61) holds, and so does the result of Lemma 29 with  $(f, \tilde{f}) = (\gamma_{\mathcal{C}}, \text{GAU}_{\mathcal{C}}(\cdot; \varepsilon))$  and  $(\nu_1, \nu_2, \mathbf{w})$  as in Algorithm 3. We will now use this fact to show (17).

By Lemma 42 (points (a) and (b)), we have  $\partial\gamma_{\mathcal{C}}(\mathbf{x}) \subseteq \mathcal{C}^\circ$  for any  $\mathbf{x} \in \mathbb{R}^d$ , and so we get

$$\sup_{\mathbf{x} \in \mathbb{R}^d} \|\partial\gamma_{\mathcal{C}}(\mathbf{x})\| \leq \sup_{\mathbf{y} \in \mathcal{C}^\circ} \|\mathbf{y}\| \leq \frac{1}{r},$$

where the last inequality follows by Lemma 42-(c). This implies that the Lipschitz constant  $L$  in Lemma 29 can be set to  $1/r$ . Furthermore, by our choice of  $\varepsilon$  and  $\nu_1$  in Algorithm 3, we have  $\varepsilon = \nu_1^3/(R^2r\sqrt{d})$  and so Lemma 29 implies that there exists a non-negative random variable  $X$  satisfying  $\mathbb{E}[X] \leq 3\nu_1 d/(rR)$ , and for all  $\mathbf{u} \in \mathbb{R}^d$

$$\gamma_{\mathcal{C}}(\mathbf{u}) \geq \gamma_{\mathcal{C}}(\mathbf{w}) + \langle \tilde{\mathbf{s}}, \mathbf{u} - \mathbf{w} \rangle - X\|\mathbf{w} - \mathbf{u}\|_\infty - 4d\nu_1/r,$$

where  $\tilde{\mathbf{s}}$  is as in Algorithm 3. Thus, with  $\Delta := 2RX + 4d\nu_1/r \geq 0$ , we get

$$\gamma_{\mathcal{C}}(\mathbf{u}) \geq \gamma_{\mathcal{C}}(\mathbf{w}) + \langle \tilde{\mathbf{s}}, \mathbf{u} - \mathbf{w} \rangle - \Delta \cdot \max(1, \|\mathbf{u}\|/R), \quad \text{and} \quad \mathbb{E}[\Delta] \leq 10d\nu_1/r, \quad (63)$$

where we used the fact that  $\mathbf{w} \in \mathcal{B}(R)$ . Further, Lemma 29 implies  $\|\tilde{\mathbf{s}}\|_\infty \leq 1/r + d^{1/4}\sqrt{L\varepsilon/\nu_1}$ ,  $\|\tilde{\mathbf{s}}\| \leq X + 1/r$ , and  $\|\tilde{\mathbf{s}}\|^2 \leq (4/r + d^{1/4}\sqrt{L\varepsilon/\nu_1})X + 1/r^2$ , and so

$$\|\tilde{\mathbf{s}}\|_\infty \leq d^{5/4}\sqrt{L\varepsilon/\nu_1} + 1/r; \quad \|\tilde{\mathbf{s}}\| \leq \Delta/R + 1/r; \quad (64)$$

$$\text{and} \quad \|\tilde{\mathbf{s}}\|^2 \leq (2/r + d^{1/4}\sqrt{L\varepsilon/\nu_1})\Delta/R + 1/r^2. \quad (65)$$

Now, since  $\varepsilon = \nu_1^3/(R^2r\sqrt{d})$  and  $\nu_1 = \delta r/(10d)$ , we get that  $10d\nu_1/r = \delta$  and  $d^{1/4}\sqrt{L\varepsilon/\nu_1} = \delta/(10Rd)$ . Plugging these into (63) and (65), we get (17) for the case when  $\tilde{\gamma} \geq 1$ .

It remain to show that  $\|\tilde{\mathbf{s}}\|_\infty < +\infty$  almost surely and bound  $\Delta$ . We do not assume that  $\tilde{\gamma} \geq 1$  anymore. Let  $i \in [n]$  and  $(\mathbf{w}_i, \mathbf{w}'_i)$  be as in Algorithm 3. Let  $\mathbf{v} \in \{\mathbf{w}_i, \mathbf{w}'_i\}$  and suppose that  $\text{GAU}_{\mathcal{C}}(\mathbf{w}; \varepsilon) \geq 1$ , then by Lemma 10 and Lemma 42-(c), we have

$$1 \leq \text{GAU}_{\mathcal{C}}(\mathbf{v}; \varepsilon) \leq \gamma_{\mathcal{C}}(\mathbf{v}) + \varepsilon \leq \|\mathbf{v}\|/r + \varepsilon \leq \frac{R + \sqrt{d}\nu_1 + \sqrt{d}\nu_2}{r} + \varepsilon < +\infty.$$

Alternatively, we have  $0 \leq \text{GAU}_{\mathcal{C}}(\mathbf{v}; \varepsilon) < 1$ . Therefore, for all  $i \in [d]$ , we have

$$\begin{aligned} |\tilde{s}_i| &= \frac{|\text{GAU}_{\mathcal{C}}(\mathbf{w}'_i; \varepsilon) - \text{GAU}_{\mathcal{C}}(\mathbf{w}_i; \varepsilon)|}{2\nu_2} \leq \frac{R + \sqrt{d}\nu_1 + \sqrt{d}\nu_2}{\nu_2 r} + \frac{\varepsilon}{\nu_2} \\ &= \frac{100d^{5/2}\kappa^2}{\delta^2 r} + \frac{10d^2\kappa}{\delta r} + \frac{\delta}{10dR} + \frac{\sqrt{d}}{r}, \\ &\leq \frac{112d^{5/2}\kappa^2}{r\delta^2} < +\infty. \end{aligned} \quad (66)$$

Finally, by Lemma 29, we have  $X \leq \|\tilde{\mathbf{s}}\|_1 + d/r$  and so combining this with (66), we get

$$\Delta = 2RX + 4d\nu_1/r \leq \frac{224d^4\kappa^3}{\delta^2} + 2d\kappa + \frac{4\delta}{10} \leq \frac{15^2d^4\kappa^3}{\delta^2},$$

where the last inequality follows by the fact that  $\delta \leq 1/3$ . This completes the proof.  $\blacksquare$

### H.3. Proof of Lemma 13 (Efficient Stochastic LO Oracle on $\mathcal{C}^\circ$ )

**Proof of Lemma 13.** Let  $I \in [d]$  be the random variable in Algorithm 4 generated during the call to  $\text{OPT}_{1d, \mathcal{C}^\circ}$  in the lemma's statement. Note that the approximate Gauge function  $\text{GAU}_{\mathcal{C}}$  (Alg. 2) is deterministic and so  $\hat{\gamma} = \tilde{\gamma}$ . On the other hand, by definition of  $\tilde{\mathbf{s}}$  and  $\hat{\mathbf{s}}$  in Algorithms 3 and 4, respectively, we have

$$\mathbb{E}[\hat{\mathbf{s}}] = \mathbb{E}[\hat{\mathbf{s}}_I \cdot \mathbf{e}_I] = \sum_{i=1}^d d^{-1} \mathbb{E}[\hat{\mathbf{s}}_i | I = i] \cdot \mathbf{e}_i \stackrel{(*)}{=} \sum_{i=1}^d d^{-1} \mathbb{E}[d\tilde{\mathbf{s}}_i] \cdot \mathbf{e}_i = \tilde{\mathbf{s}},$$

where  $(*)$  follows by the fact that, conditioned on  $I = i$ ,  $\hat{\mathbf{s}}_i/d$  has the same distribution as  $\tilde{\mathbf{s}}_i$ . Similarly,

$$\mathbb{E}[\|\hat{\mathbf{s}}\|^2] = \mathbb{E}[\|\hat{\mathbf{s}}_I \mathbf{e}_I\|^2] = \sum_{i=1}^d d^{-1} \mathbb{E}[\hat{\mathbf{s}}_i^2 \mathbf{e}_i | I = i] = \sum_{i=1}^d d^{-1} \mathbb{E}[d^2 \tilde{\mathbf{s}}_i^2 \mathbf{e}_i] = d \mathbb{E}[\|\tilde{\mathbf{s}}\|^2] \leq d \cdot \left( \frac{1}{r} + \frac{\delta}{R} \right)^2,$$

where the last inequality follows by the fact that  $\|\tilde{\mathbf{s}}\|^2 \leq (2/r + \delta/R)\Delta/R + 1/r^2$  (Proposition 11), where  $\Delta \geq 0$  is a random variable satisfying  $\mathbb{E}[\Delta] \leq \delta$ . Finally, the fact that  $\|\hat{\mathbf{s}}\| < +\infty$  a.s. follows from the fact that  $\|\tilde{\mathbf{s}}\|_\infty < +\infty$  a.s. (Prop. 11), and for any  $i \in [d]$ , conditioned on  $I = i$ ,  $\hat{\mathbf{s}}_i/d$  has the same distribution as  $\tilde{\mathbf{s}}_i$ .  $\blacksquare$

## Appendix I. Proofs of the Regret Bounds and Convergence Rates

### I.1. Proof of Lemma 7 (Instantaneous Regret Bound)

To avoid expensive projections our main algorithm (Alg. 1) makes use of surrogate losses of the form  $\tilde{\ell}(\mathbf{w}) := \langle \mathbf{g}, \mathbf{w} \rangle + bS_{\mathcal{C}}(\mathbf{w})$ ,  $\mathbf{w} \in \mathbb{R}^d$  and  $b \geq 0$ . The choice of such a surrogate loss function is inspired by existing constrained-to-unconstrained reductions in OCO due to Cutkosky and Orabona (2018); Cutkosky (2020). We will use the approximate optimization Oracle  $\text{OPT}_{\mathcal{C}^\circ}$  from App. B to compute approximate subgradients of such surrogate losses. In particular,  $\text{OPT}_{\mathcal{C}^\circ}$  (Algorithm 3) guarantees the following:

**Lemma 30** *Let  $\mathbf{g} \in \mathbb{R}^d$ ,  $b \geq 0$ , and  $\delta \in (0, 1/3)$ . Let  $S_{\mathcal{C}}$  be the Gauge distance function in (3) and define  $\tilde{\ell}(\mathbf{w}) := \langle \mathbf{g}, \mathbf{w} \rangle + bS_{\mathcal{C}}(\mathbf{w})$ . Then, for  $\mathbf{w} \in \mathcal{B}(R)$ , the output  $(\tilde{\gamma}, \tilde{\mathbf{s}})$  of Alg. 3 with input  $(\mathbf{w}, \delta)$  satisfies (17) and*

$$\forall \mathbf{u} \in \mathcal{B}(R), \quad \tilde{\ell}(\mathbf{u}) \geq \tilde{\ell}(\mathbf{w}) + \langle \mathbf{g} + b\mathbf{v}, \mathbf{u} - \mathbf{w} \rangle - b \cdot (\Delta + \delta) \cdot \mathbb{I}_{\{\tilde{\gamma} \geq 1\}}, \quad \text{where } \mathbf{v} := \mathbb{I}_{\{\tilde{\gamma} \geq 1\}} \tilde{\mathbf{s}},$$

and  $\Delta \in [0, 15^2d^4\kappa^3\delta^{-2}]$  is the same random variable satisfying (17); in particular,  $\mathbb{E}[\Delta] \leq \delta$ .

**Proof** First suppose that  $\tilde{\gamma} < 1$ . If  $\gamma_{\mathcal{C}}(\mathbf{w}) \geq 9/16$ , then by Lemma 10, we have  $\gamma_{\mathcal{C}}(\mathbf{w}) \leq \tilde{\gamma} < 1$ . Alternatively,  $\gamma_{\mathcal{C}}(\mathbf{w}) < 9/16 < 1$ . Therefore, by Lemma 6,  $\partial S_{\mathcal{C}}(\mathbf{w}) = \{\mathbf{0}\}$ , and so  $\mathbf{g} \in \partial \tilde{\ell}(\mathbf{w})$ . Since the function  $\tilde{\ell}$  is convex, it follows that

$$[\text{case where } \tilde{\gamma} < 1] \quad \forall \mathbf{u} \in \mathcal{B}(R), \quad \tilde{\ell}(\mathbf{u}) \geq \tilde{\ell}(\mathbf{w}) + \langle \mathbf{g}, \mathbf{u} - \mathbf{w} \rangle. \quad (67)$$

Now, suppose that  $\tilde{\gamma} \geq 1$ . In this case, by Lemma 10, we have  $\gamma_{\mathcal{C}}(\mathbf{w}) \leq \tilde{\gamma} \leq \gamma_{\mathcal{C}}(\mathbf{w}) + \delta$ . Combining this with the fact that  $S_{\mathcal{C}}(\mathbf{w}) = 0 \vee (\gamma_{\mathcal{C}}(\mathbf{w}) - 1)$  (Lemma 6), we get

$$\gamma_{\mathcal{C}}(\mathbf{w}) - 1 \geq \tilde{\gamma} - 1 - \delta \stackrel{(a)}{\geq} 0 \vee (\gamma_{\mathcal{C}}(\mathbf{w}) - 1) - \delta = S_{\mathcal{C}}(\mathbf{w}) - \delta, \quad (68)$$

where (a) follows from the fact that  $\tilde{\gamma} \geq 1$  (by assumption) and  $\tilde{\gamma} \geq \gamma_{\mathcal{C}}(\mathbf{w})$ . Using this, we get

$$\begin{aligned} [\text{case where } \tilde{\gamma} \geq 1] \quad \forall \mathbf{u} \in \mathcal{B}(R), \quad \tilde{\ell}(\mathbf{u}) &= \langle \mathbf{g}, \mathbf{u} \rangle + 0 \vee (b\gamma_{\mathcal{C}}(\mathbf{u}) - b), \\ &\geq \langle \mathbf{g}, \mathbf{u} \rangle + b\gamma_{\mathcal{C}}(\mathbf{u}) - b, \\ &\geq \langle \mathbf{g}, \mathbf{u} \rangle + b\gamma_{\mathcal{C}}(\mathbf{w}) - b + \langle b\tilde{\mathbf{s}}, \mathbf{u} - \mathbf{w} \rangle - b\Delta, \end{aligned} \quad (69)$$

$$\geq \langle \mathbf{g}, \mathbf{u} \rangle + bS_{\mathcal{C}}(\mathbf{w}) + \langle b\tilde{\mathbf{s}}, \mathbf{u} - \mathbf{w} \rangle - b\Delta - b\delta, \quad (70)$$

$$\begin{aligned} &= \tilde{\ell}(\mathbf{w}) + \langle \mathbf{g} + b\tilde{\mathbf{s}}, \mathbf{u} - \mathbf{w} \rangle - b\Delta - b\delta, \\ &= \tilde{\ell}(\mathbf{w}) + \langle \tilde{\mathbf{g}} + b\mathbf{v}, \mathbf{u} - \mathbf{w} \rangle - b\Delta - b\delta, \end{aligned} \quad (71)$$

where (69) follows by the fact that  $\tilde{\gamma} \geq 1$  and Proposition 11; (70) follows by (68); and (71) follows by the fact that  $\tilde{\mathbf{s}} = \mathbf{v}$  since  $\tilde{\gamma} \geq 1$ . Combining (67) and (71) implies the desired result.  $\blacksquare$

Lemma 30 shows that for any  $t \geq 1$ , the vector  $\tilde{\mathbf{g}}_t$  in Algorithm 1 is an approximate subgradient of the surrogate loss  $\tilde{\ell}_t(\cdot) := \langle \mathbf{g}_t, \cdot \rangle - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle S_{\mathcal{C}}(\cdot)$ , where  $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$  and  $\mathbf{x}_t$  is the  $t$ th iterate of Alg. 1. We now use this to prove Lemma 7. We actually state and prove a slight generalization of Lemma 7, which will be useful when considering the stochastic optimization setting of Section C.3. In this generalization, we assume that for some  $R' \in [r, R]$ ,

$$\mathcal{B}(r) \subseteq \mathcal{C} \subseteq \mathcal{B}(R') \subseteq \mathcal{B}(R). \quad (72)$$

**Lemma 31** *Let  $\kappa := R'/r$  and  $\mathbf{A}$  be any OCO algorithm on  $\mathcal{B}(R)$ , for  $r, R, R' > 0$  as in (72). Further, for  $t \geq 1$ , let  $\mathbf{w}_t, \tilde{\mathbf{g}}_t, \mathbf{x}_t$ , and  $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$  be as in Alg. 1 with  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{\mathcal{C}^\circ}$  and any tolerance sequence  $(\delta_s) \subset (0, 1/3)$ . Then,  $(\mathbf{x}_t) \subset \mathcal{C}$ , and for all  $t \geq 1$ , there exists a variable  $\Delta_t \in [0, 15^2 d^4 (R/r)^3 \delta_t^{-2}]$  s.t.  $\mathbb{E}_{t-1}[\Delta_t] \leq \delta_t$ ,  $\|\tilde{\mathbf{g}}_t\| \leq (1 + \Delta_t + \kappa)\|\mathbf{g}_t\|$ , and*

$$\forall \mathbf{x} \in \mathcal{C}, \quad \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \tilde{\ell}_t(\mathbf{w}_t) - \tilde{\ell}_t(\mathbf{x}) + \delta_t R \|\mathbf{g}_t\| \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + (2\delta_t + \Delta_t) R \|\mathbf{g}_t\|. \quad (73)$$

**Proof of Lemma 7.** Let  $\gamma_t, \mathbf{v}_t, \mathbf{w}_t, \tilde{\mathbf{g}}_t$ , and  $\mathbf{x}_t$  be as in Algorithm 1. We will first show that

$$-R\|\mathbf{g}_t\| \leq \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \cdot \langle \mathbf{g}_t, \mathbf{x}_t \rangle \leq 0 \quad (74)$$

First suppose that  $\gamma_t < 1$ . In this case,  $\mathbf{x}_t = \mathbf{w}_t$ , and so

$$-R\|\mathbf{g}_t\| \leq \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \cdot \langle \mathbf{g}_t, \mathbf{x}_t \rangle = \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \cdot \langle \mathbf{g}_t, \mathbf{w}_t \rangle \leq 0.$$



Now suppose that  $\gamma_t \geq 1$ , then  $\mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \cdot \langle \mathbf{g}_t, \mathbf{x}_t \rangle = \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \cdot \langle \mathbf{g}_t, \mathbf{w}_t / \gamma_t \rangle$  and so

$$0 \geq \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \cdot \langle \mathbf{g}_t, \mathbf{x}_t \rangle \geq \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{w}_t \rangle \geq -R \|\mathbf{g}_t\|,$$

where the last inequality follows by the fact  $\mathbf{w}_t \in \mathcal{B}(R)$ . This shows (74). Now recall the definition of the surrogate function  $\tilde{\ell}_t : \mathcal{B}(R) \rightarrow \mathbb{R}$ :

$$\tilde{\ell}_t(\mathbf{w}) := \langle \mathbf{g}_t, \mathbf{w} \rangle - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \cdot \langle \mathbf{g}_t, \mathbf{x}_t \rangle \cdot S_{\mathcal{C}}(\mathbf{w}),$$

where  $S_{\mathcal{C}}$  is as in (3). By (74) and Lemma 30, there exists a r.v.  $\Delta_t \in [0, \frac{15^2 d^4 \kappa^3}{\delta^2}]$  such that  $\mathbb{E}_{t-1}[\Delta_t] \leq \delta_t$  and

$$\forall \mathbf{x} \in \mathcal{C}, \quad \tilde{\ell}_t(\mathbf{w}_t) - \tilde{\ell}_t(\mathbf{x}) \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + R \|\mathbf{g}_t\| \cdot (\Delta_t + \delta_t). \quad (75)$$

It remains to show that  $\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \tilde{\ell}_t(\mathbf{w}_t) - \tilde{\ell}_t(\mathbf{x}) + \delta_t R \|\mathbf{g}_t\|$ , for all  $t \geq 1$  and  $\mathbf{x} \in \mathcal{C}$ . First, note that for all  $\mathbf{x} \in \mathcal{C}$ , we have  $S_{\mathcal{C}}(\mathbf{x}) = 0$ , and so

$$\tilde{\ell}_t(\mathbf{x}) = \langle \mathbf{g}_t, \mathbf{x} \rangle, \quad \forall \mathbf{x} \in \mathcal{C}. \quad (76)$$

We will now compare  $\langle \mathbf{g}_t, \mathbf{x}_t \rangle$  to  $\tilde{\ell}_t(\mathbf{w}_t)$  by considering cases. Suppose that  $\gamma_t < 1$ . In this case, we have  $\mathbf{x}_t = \mathbf{w}_t$  and so

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle = \langle \mathbf{g}_t, \mathbf{w}_t \rangle = \tilde{\ell}_t(\mathbf{w}_t). \quad [\text{case where } \gamma_t < 1] \quad (77)$$

Now suppose that  $\gamma_t \geq 1$  and  $\langle \mathbf{g}_t, \mathbf{w}_t \rangle \geq 0$ . In this case, we have  $\mathbf{x}_t = \mathbf{w}_t / \gamma_t$ , and so

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle = \langle \mathbf{g}_t, \mathbf{w}_t / \gamma_t \rangle \leq \langle \mathbf{g}_t, \mathbf{w}_t \rangle = \tilde{\ell}_t(\mathbf{w}_t). \quad [\text{case where } \gamma_t \geq 1, \langle \mathbf{g}_t, \mathbf{w}_t \rangle \geq 0] \quad (78)$$

Now suppose that  $\gamma_t \geq 1$ ,  $\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0$ , and  $\mathbf{w}_t \in \mathcal{C}$ . We note that this implies that  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \leq 1$  and so  $S_{\mathcal{C}}(\mathbf{w}_t) = 0$  (Lemma 6). Thus,  $\tilde{\ell}_t(\mathbf{w}_t) = \langle \mathbf{g}_t, \mathbf{w}_t \rangle$ . On the other hand, by Lemma 10 we have  $\gamma_t \leq \gamma_{\mathcal{C}}(\mathbf{w}_t) + \delta_t \leq 1 + \delta_t$ , and so since  $\mathbf{x}_t = \mathbf{w}_t / \gamma_t$ , we have

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle \leq \langle \mathbf{g}_t, \mathbf{w}_t \rangle / (1 + \delta_t) = \langle \mathbf{g}_t, \mathbf{w}_t \rangle - \delta_t \langle \mathbf{g}_t, \mathbf{w}_t \rangle / (1 + \delta_t).$$

Thus, since  $\tilde{\ell}_t(\mathbf{w}_t) = \langle \mathbf{g}_t, \mathbf{w}_t \rangle$ , the previous display implies that

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle \leq \tilde{\ell}_t(\mathbf{w}_t) + \delta_t R \|\mathbf{g}_t\|. \quad [\text{case where } \gamma_t \geq 1, \langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0, \mathbf{w}_t \in \mathcal{C}] \quad (79)$$

We now consider the last case where  $\gamma_t \geq 1$ ,  $\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0$ , and  $\mathbf{w}_t \notin \mathcal{C}$ . We note that this implies that  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \geq 1$ . By the fact that  $\gamma_t \leq \gamma_{\mathcal{C}}(\mathbf{w}_t) + \delta_t$  (Lemma 10), we have

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle \leq \frac{\langle \mathbf{g}_t, \mathbf{w}_t \rangle}{\gamma_{\mathcal{C}}(\mathbf{w}_t) + \delta_t} = \frac{\langle \mathbf{g}_t, \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t) \rangle}{1 + \delta_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)} = \langle \mathbf{g}_t, \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t) \rangle - \frac{\delta_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)}{1 + \delta_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)} \langle \mathbf{g}_t, \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t) \rangle, \quad (80)$$

Thus, since  $S_{\mathcal{C}}(\mathbf{w}_t) = \gamma_{\mathcal{C}}(\mathbf{w}_t) - 1$  (by Lemma 6 and  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \geq 1$ ), we get

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle \cdot S_{\mathcal{C}}(\mathbf{w}_t) \leq \langle \mathbf{g}_t, \mathbf{w}_t \rangle - \frac{\delta_t \langle \mathbf{g}_t, \mathbf{w}_t \rangle}{1 + \delta_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)} - \langle \mathbf{g}_t, \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t) \rangle + \frac{\delta_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)}{1 + \delta_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)} \langle \mathbf{g}_t, \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t) \rangle.$$

Adding this together with (80) and using the fact that  $\|\mathbf{w}_t\| \leq R$ , we get

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle + \langle \mathbf{g}_t, \mathbf{x}_t \rangle S_{\mathcal{C}}(\mathbf{w}_t) \leq \langle \mathbf{g}_t, \mathbf{w}_t \rangle + \delta_t R \|\mathbf{g}_t\|,$$

and so after rearranging and using that  $\tilde{\ell}_t(\mathbf{w}_t) = \langle \mathbf{g}_t, \mathbf{w}_t \rangle - \langle \mathbf{g}_t, \mathbf{x}_t \rangle S_{\mathcal{C}}(\mathbf{w}_t)$ , we get

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle \leq \tilde{\ell}_t(\mathbf{w}_t) + \delta_t R \|\mathbf{g}_t\|. \quad [\text{case where } \gamma_t \geq 1, \langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0, \mathbf{w}_t \notin \mathcal{C}] \quad (81)$$

By combining, (77), (78), (79), and (81), we obtain:

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle \leq \tilde{\ell}_t(\mathbf{w}_t) + \delta_t R \|\mathbf{g}_t\|.$$

Combining this with (75) and (76), we get

$$\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \tilde{\ell}_t(\mathbf{w}_t) - \tilde{\ell}_t(\mathbf{x}) + \delta_t R \|\mathbf{g}_t\| \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + (2\delta_t + \Delta_t) R \|\mathbf{g}_t\|, \quad \forall \mathbf{x} \in \mathcal{C}.$$

This shows (73). It remains to bound  $\|\tilde{\mathbf{g}}_t\|$  in terms of  $\|\mathbf{g}_t\|$  and show that  $\mathbf{x}_t \in \mathcal{C}$ . When  $\gamma_t < 1$ , we have  $\tilde{\mathbf{g}}_t = \mathbf{g}_t$  and  $\mathbf{x}_t = \mathbf{w}_t$ , and so  $\|\tilde{\mathbf{g}}_t\| = \|\mathbf{g}_t\|$ . Furthermore, we also have that  $\mathbf{x}_t \in \mathcal{C}$ . In fact, if  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \leq 9/16$ , then by definition of the Gauge function we have  $\mathbf{w}_t \in \mathcal{C}$  and the same holds for  $\mathbf{x}_t$  (since  $\mathbf{x}_t = \mathbf{w}_t$ ). On the other hand, if  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \geq 9/16$ , then by Lemma 10, we have  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \leq \gamma_t < 1$  ( $\gamma_t < 1$  is the case we are currently considering), and so  $\mathbf{x}_t = \mathbf{w}_t \in \mathcal{C}$ .

Now suppose that  $\gamma_t \geq 1$ . In this case, we have  $\mathbf{x}_t = \mathbf{w}_t/\gamma_t$  and  $\tilde{\mathbf{g}}_t = \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \mathbf{v}_t$ . The latter fact together with the positive homogeneity of the Gauge function (Lemma 42-(a,b)) imply that  $\gamma_{\mathcal{C}}(\mathbf{x}_t) = \gamma_{\mathcal{C}}(\mathbf{w}_t/\gamma_t) = \gamma_{\mathcal{C}}(\mathbf{w}_t)/\gamma_t \leq \gamma_{\mathcal{C}}(\mathbf{w}_t)/\gamma_{\mathcal{C}}(\mathbf{w}_t) = 1$  (since  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \leq \gamma_t$  by Lemma 10), and so  $\mathbf{x}_t \in \mathcal{C}$ . Using that  $\mathbf{x}_t \in \mathcal{C}$  (which implies that  $\|\mathbf{x}_t\| \leq R'$ ) and that  $\tilde{\mathbf{g}}_t = \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \mathbf{v}_t$ , we get

$$\|\tilde{\mathbf{g}}_t\| = \|\mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \mathbf{v}_t\| \leq \|\mathbf{g}_t\| (1 + R' \|\mathbf{v}_t\|) \leq (1 + \Delta_t + \kappa) \|\mathbf{g}_t\|,$$

where the last inequality follows the fact that  $\|\mathbf{v}_t\| \leq \Delta_t/R + 1/r$  (by (17) which is implied by Lemma 30). ■

**Proof of Lemma 7.** Follows from Lemma 31 with  $R' = R$ . ■

## I.2. Proof of Theorem 8 (Regret Bound in High Probability using $\text{OPT}_{\mathcal{C}^\circ}$ )

**Proof** By Lemma 7, we have, for all  $\mathbf{w}$  and  $t \geq 1$ ,  $\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + (2\delta_t + \Delta_t) R \|\mathbf{g}_t\|$ , where  $(\Delta_t) \subset \mathbb{R}_{\geq 0}$  is a sequence of positive random variables satisfying  $\mathbb{E}_{t-1}[\Delta_t] \leq \delta_t$ , for all  $t \in [T]$ . Summing this inequality for  $t = 1, \dots, T$ , we obtain, for all  $\mathbf{x} \in \mathcal{C}$ ,

$$\begin{aligned} \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle &\leq \sum_{t=1}^T \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + \sum_{t=1}^T (2\delta_t + \Delta_t) R \|\mathbf{g}_t\|, \\ &\leq 2R \sqrt{2 \sum_{t=1}^T \|\tilde{\mathbf{g}}_t\|^2 + \sum_{t=1}^T (2\delta_t + \Delta_t) R \|\mathbf{g}_t\|}, \end{aligned} \quad (82)$$

$$\leq 4(1 + \kappa) R \sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2 + \sum_{t=1}^T (2\delta_t + 5\Delta_t) R \|\mathbf{g}_t\|}, \quad (83)$$

where (82) follows by our choice of the subroutine A and the regret bound of FTRL-prox in Proposition 25, and the last inequality follows by the fact that  $\|\tilde{\mathbf{g}}_t\| \leq (1 + \kappa + \Delta_t)\|\mathbf{g}_t\|$  (Lem. 7) and the fact that  $(a + b)^2 \leq 2a^2 + 2b^2$  for all  $a, b \in \mathbb{R}_{>0}$ . By Lemma 7, we also have  $\mathbf{x}_t \in \mathcal{C}$ ,  $\forall t \geq 1$ .

We now instantiate (83) with the specific choice of tolerance sequence  $(\delta_t)$  in theorem's statement, where we explicitly bound the right-most sum in (82) involving  $(\Delta_t)$ . For this, we let  $X_t := \sum_{i=1}^t (\Delta_i - \bar{\delta}_i)$ , where  $\bar{\delta}_i := \mathbb{E}_{i-1}[\Delta_i] \leq \delta_i$ . The process  $(X_t)$  is a martingale; that is, for all  $i \geq 1$ , we have,  $\mathbb{E}_i[X_t] = X_i$ , for all  $i < t$ . Thus, by Doob's martingale inequality (Durrett, 2019, Theorem 4.4.2) we have, for any  $\rho \in (0, 1)$  and  $T \geq 1$ :

$$\begin{aligned} \mathbb{P}\left[\sum_{t=1}^T \Delta_t \geq (1 + 1/\rho) \sum_{t=1}^T \delta_t\right] &\leq \mathbb{P}\left[X_T \geq \sum_{t=1}^T \delta_t/\rho\right] \leq \mathbb{P}\left[\max_{t \leq T} X_t \geq \sum_{t=1}^T \delta_t/\rho\right] \\ &\leq \frac{\rho \mathbb{E}[X_T \vee 0]}{\sum_{t=1}^T \delta_t} \leq \rho, \end{aligned}$$

where the last inequality follows by the fact that  $\mathbb{E}[X_T \vee 0] \leq \mathbb{E}[\sum_{t=1}^T \Delta_t] \leq \sum_{t=1}^T \delta_t$ . Using this and the fact that  $\sum_{t=1}^{\infty} 1/t^2 \leq 2$  in combination with (83), we obtain (9).  $\blacksquare$

### I.3. Proof of Theorem 15 (Regret Bound in Expectation using $\text{OPT}_{1d, \mathcal{C}^\circ}$ )

**Proof** Let  $(\tilde{\gamma}_t, \tilde{\mathbf{s}}_t) = \text{OPT}_{\mathcal{C}^\circ}(\mathbf{w}_t; \delta_t)$  and  $(\hat{\gamma}_t, \hat{\mathbf{s}}_t) = \text{OPT}_{1d, \mathcal{C}^\circ}(\mathbf{w}_t; \delta_t)$ . Further, let  $\tilde{\mathbf{v}}_t := \mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \tilde{\mathbf{s}}_t$ , and  $\gamma_t$  and  $\mathbf{v}_t$  be as in Algorithm 1. Note that by Proposition 11 and Lemma 13, we have  $\|\tilde{\mathbf{s}}_t\| \vee \|\hat{\mathbf{s}}_t\| < +\infty$  almost surely and so the expectations  $\mathbb{E}_{t-1}[\mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \tilde{\mathbf{s}}_t]$  and  $\mathbb{E}_{t-1}[\mathbb{I}_{\{\hat{\gamma}_t \geq 1\}} \hat{\mathbf{s}}_t]$  are well defined. We also note that  $\gamma_t = \hat{\gamma}_t = \tilde{\gamma}_t$ , and  $\gamma_t, \mathbf{w}_t, \mathbf{g}_t$ , and  $\mathbf{x}_t$  are all deterministic functions of the past ( $\mathbf{w}_t$  is the output of FTRL-prox, which is a deterministic function of the past).

By Lemma 7, there exists a random variable  $\Delta_t \geq 0$  satisfying  $\mathbb{E}_{t-1}[\Delta_t]$  and such that for all  $t \geq 1$ ,

$$\forall \mathbf{x} \in \mathcal{C}, \quad \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \langle \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \tilde{\mathbf{v}}_t, \mathbf{w}_t - \mathbf{x} \rangle + (2\delta_t + \Delta_t)R\|\mathbf{g}_t\|, \quad (84)$$

Now by Lemma 13, and the law of total expectation, we have, for all  $\mathbf{x} \in \mathcal{C}$ ,

$$\begin{aligned} \mathbb{E}[\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle] &= \mathbb{E}[\mathbb{E}_{t-1}[\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle]], \\ &\leq \mathbb{E}[\mathbb{E}_{t-1}[\langle \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \tilde{\mathbf{v}}_t, \mathbf{w}_t - \mathbf{x} \rangle + (2\delta_t + \Delta_t)R\|\mathbf{g}_t\|]], \quad (\text{by (84)}) \\ &\leq \mathbb{E}[\mathbb{E}_{t-1}[\langle \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \tilde{\mathbf{v}}_t, \mathbf{w}_t - \mathbf{x} \rangle] + 3\delta_t R\|\mathbf{g}_t\|], \\ &= \mathbb{E}[\mathbb{I}_{\{\tilde{\gamma}_t < 1\}} \cdot \mathbb{E}_{t-1}[\langle \mathbf{g}_t, \mathbf{w}_t - \mathbf{x} \rangle]] \\ &\quad + \mathbb{E}[\mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \cdot \mathbb{E}_{t-1}[\langle \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \tilde{\mathbf{s}}_t, \mathbf{w}_t - \mathbf{x} \rangle] + 3\delta_t R\|\mathbf{g}_t\|], \\ &= \mathbb{E}[\mathbb{I}_{\{\tilde{\gamma}_t < 1\}} \cdot \mathbb{E}_{t-1}[\langle \mathbf{g}_t, \mathbf{w}_t - \mathbf{x} \rangle]] \\ &\quad + \mathbb{E}[\mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \cdot \mathbb{E}_{t-1}[\langle \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \hat{\mathbf{s}}_t, \mathbf{w}_t - \mathbf{x} \rangle] + 3\delta_t R\|\mathbf{g}_t\|], \quad (85) \\ &= \mathbb{E}[\mathbb{E}_{t-1}[\langle \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \mathbf{v}_t, \mathbf{w}_t - \mathbf{x} \rangle] + 3\delta_t R\|\mathbf{g}_t\|], \\ &= \mathbb{E}[\langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + 3\delta_t R\|\mathbf{g}_t\|], \quad (86) \end{aligned}$$

where (85) follows by the facts that  $\tilde{\gamma}_t = \hat{\gamma}_t$  and  $\mathbb{E}_{t-1}[\hat{s}_t] = \mathbb{E}_{t-1}[\tilde{s}_t]$ , when  $\tilde{\gamma}_t \geq 1$  (Lemma 13). Summing (86) for  $t = 1, \dots, T$ , we obtain,

$$\begin{aligned} \forall \mathbf{x} \in \mathcal{C}, \quad \sum_{t=1}^T \mathbb{E}[\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle] &\leq \mathbb{E} \left[ \sum_{t=1}^T \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + 3 \sum_{t=1}^T \delta_t R \|\mathbf{g}_t\| \right], \\ &\leq \mathbb{E} \left[ 2R \sqrt{2 \sum_{t=1}^T \|\tilde{\mathbf{g}}_t\|^2} + 3 \sum_{t=1}^T \delta_t R \|\mathbf{g}_t\| \right], \end{aligned} \quad (87)$$

$$\leq 2R \sqrt{\mathbb{E} \left[ 2 \sum_{t=1}^T \mathbb{E}_{t-1} [\|\tilde{\mathbf{g}}_t\|^2] \right]} + \mathbb{E} \left[ 3 \sum_{t=1}^T \delta_t R \|\mathbf{g}_t\| \right], \quad (88)$$

where the last step follows by Jensen's inequality and (87) follows by our choice of the subroutine A ( $\equiv$ FTRL-prox) and the regret bound of FTRL-prox in Proposition 25.

It remains to bound  $\|\tilde{\mathbf{g}}_t\|$  in terms of  $\|\mathbf{g}_t\|$  and show that  $\mathbf{x}_t \in \mathcal{C}$ . First, if  $\gamma_t < 1$ , then  $(\tilde{\mathbf{g}}_t, \mathbf{x}_t) = (\mathbf{g}_t, \mathbf{w}_t)$ , and so  $\|\tilde{\mathbf{g}}_t\| \leq \|\mathbf{g}_t\|(1 + R\|\mathbf{s}_t\|)$ . Furthermore, we also have that  $\mathbf{x}_t \in \mathcal{C}$ . In fact, if  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \leq 9/16$ , then by definition of the Gauge function we have  $\mathbf{w}_t \in \mathcal{C}$  and the same holds for  $\mathbf{x}_t$  (since  $\mathbf{x}_t = \mathbf{w}_t$ ). On the other hand, if  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \geq 9/16$ , then by Lemma 10, we have  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \leq \gamma_t < 1$ , and so  $\mathbf{x}_t = \mathbf{w}_t \in \mathcal{C}$ .

Now suppose that  $\gamma_t \geq 1$ . In this case, we have  $\mathbf{x}_t = \mathbf{w}_t/\gamma_t$  and  $\mathbf{v}_t = \mathbf{s}_t$ . Using this, the triangular inequality, and Cauchy Schwarz, we get

$$\|\tilde{\mathbf{g}}_t\| = \|\langle \mathbf{g}_t, \mathbf{w} \rangle - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \cdot \langle \mathbf{g}_t, \mathbf{x}_t \rangle \cdot \mathbf{s}_t\| = \|\mathbf{g}_t\| + \|\mathbf{g}_t\| \|\mathbf{w}_t\|/\gamma_t \cdot \|\mathbf{s}_t\| \leq \|\mathbf{g}_t\|(1 + R\|\mathbf{s}_t\|),$$

where the last inequality follows by the fact that  $\gamma_t \geq 1$ . Therefore, using Lemma 13, we get  $\|\tilde{\mathbf{g}}_t\| \leq \|\mathbf{g}_t\|(1 + \Delta_t + \kappa)$  and

$$\mathbb{E}_{t-1}[\|\tilde{\mathbf{g}}_t\|^2] \leq 2\|\mathbf{g}_t\|^2(1 + R^2\mathbb{E}_{t-1}[\|\mathbf{s}_t\|^2]) = 2\|\mathbf{g}_t\|^2(1 + R^2\mathbb{E}_{t-1}[\|\hat{\mathbf{s}}_t\|^2]) \leq 2\|\mathbf{g}_t\|^2(1 + d \cdot (\kappa + \delta_t)^2).$$

Plugging this into (88) leads to, for all  $\mathbf{x} \in \mathcal{C}$ ,

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle] &\leq 4R \sqrt{\sum_{t=1}^T (1 + d \cdot (\kappa + \delta_t)^2) \cdot \mathbb{E}[\|\mathbf{g}_t\|^2]} + 3\mathbb{E} \left[ \sum_{t=1}^T \delta_t R \cdot \|\mathbf{g}_t\| \right], \\ &\leq 4R \sqrt{\sum_{t=1}^T (1 + d \cdot (\kappa + \delta)^2) \cdot \mathbb{E}[\|\mathbf{g}_t\|^2]} + 6\delta R \cdot \mathbb{E} \left[ \max_{t \in [T]} \|\mathbf{g}_t\| \right], \end{aligned}$$

where in the last inequality we used that  $\sum_{t=1}^{+\infty} 1/t^2 \leq 2$ . Furthermore, when  $\gamma_t \geq 1$ , we have  $\gamma_{\mathcal{C}}(\mathbf{x}_t) = \gamma_{\mathcal{C}}(\mathbf{w}_t/\gamma_t) = \gamma_{\mathcal{C}}(\mathbf{w}_t)/\gamma_t \leq \gamma_{\mathcal{C}}(\mathbf{w}_t)/\gamma_{\mathcal{C}}(\mathbf{w}_t) = 1$  (since  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \leq \gamma_t$  by Lemma 10), and so  $\mathbf{x}_t \in \mathcal{C}$ .  $\blacksquare$

#### I.4. Proof of Theorem 16 (Regret Bound in High Probability using $\text{OPT}_{\text{Id}, \mathcal{C}^\circ}$ )

To prove Theorem 16, we need the following extension of Lemma 13:

**Lemma 32** *Let  $\delta \in (0, 1)$ ,  $\mathbf{w} \in \mathcal{B}(R)$ , and  $\kappa := R/r$ , with  $r, R > 0$  as in (1). Further, let  $(\tilde{\gamma}, \tilde{\mathbf{s}}) = \text{OPT}_{\mathcal{C}^\circ}(\mathbf{w}; \delta)$  and  $(\hat{\gamma}, \hat{\mathbf{s}}) = \text{OPT}_{1d, \mathcal{C}^\circ}(\mathbf{w}; \delta)$  (Alg. 4). Then,  $\hat{\gamma} = \tilde{\gamma}$ ,  $\|\hat{\mathbf{s}}\| < +\infty$  a.s., and if  $\hat{\gamma} \geq 1$  it follows that*

$$\mathbb{E}[\hat{\mathbf{s}}] = \mathbb{E}[\tilde{\mathbf{s}}], \quad \mathbb{E}[\|\hat{\mathbf{s}}\|^2] \leq d \cdot (1/r + \delta/R)^2, \quad (89)$$

$$\|\hat{\mathbf{s}}\|_\infty \leq d \cdot (1/r + \delta/R), \quad \|\hat{\mathbf{s}}\| \leq d \cdot (1/r + \delta/R), \quad \text{and} \quad \mathbb{E}[\|\hat{\mathbf{s}}\|^4] \leq d^3 \cdot (1/r + \delta/R)^4. \quad (90)$$

**Proof** Lemma 13 implies (89). We now show (90). Proposition 11 implies that  $\tilde{s}_i \leq \delta/R + 1/r$ , for all  $i \in [d]$ . Let  $I \in [d]$  be the random variable in Algorithm 4 generated during the call to  $\text{OPT}_{1d, \mathcal{C}^\circ}$  in the lemma's statement. Since, conditioned on  $I = i$ ,  $\hat{s}_i/d$  has the same distribution as  $\tilde{s}_i$  and  $\hat{s}_i = 0$  when  $I \neq i$ , we get that  $\hat{s}_i \leq d \cdot (1/r + \delta/R)$ . This implies the first inequality in (90). Using this, we get

$$\begin{aligned} \mathbb{E}[\|\hat{\mathbf{s}}\|^4] &= \mathbb{E}[\|\hat{s}_I \mathbf{e}_I\|^4] = \mathbb{E}[\hat{s}_I^4] \leq d^2 \cdot \left(\frac{1}{r} + \frac{\delta}{R}\right)^2 \cdot \mathbb{E}[\hat{s}_I^2] = d^2 \cdot \left(\frac{1}{r} + \frac{\delta}{R}\right)^2 \cdot \mathbb{E}[\|\hat{\mathbf{s}}\|^2] \\ &\leq d^3 \cdot \left(\frac{1}{r} + \frac{\delta}{R}\right)^4, \end{aligned}$$

where the last inequality follows the right-most inequality in (89). This shows the right-most inequality in (90). Finally, we have

$$\|\hat{\mathbf{s}}\| = |\hat{s}_I| \leq \max_{i \in [d]} |\hat{s}_i| \leq d \cdot (1/r + \delta/R),$$

where the last inequality follows by the first inequality in (90), which we already showed. Now we do not assume that  $\hat{\gamma} \geq 1$  anymore. Finally, the fact that  $\|\hat{\mathbf{s}}\| < +\infty$  follows by Lemma 13.  $\blacksquare$

**Proof of Theorem 16.** The fact that  $(\mathbf{x}_t) \subset \mathcal{C}$  follows from Lemma 15. Let  $(\tilde{\gamma}_t, \tilde{\mathbf{s}}_t) = \text{OPT}_{\mathcal{C}^\circ}(\mathbf{w}_t; \delta_t)$  and  $(\hat{\gamma}_t, \hat{\mathbf{s}}_t) = \text{OPT}_{1d, \mathcal{C}^\circ}(\mathbf{w}_t; \delta_t)$ . Further, let  $\tilde{\mathbf{v}}_t := \mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \tilde{\mathbf{s}}_t$ , and  $\gamma_t$  and  $\mathbf{v}_t$  be as in Algorithm 1. Note that by Proposition 11 and Lemma 13, we have  $\|\tilde{\mathbf{s}}_t\| \vee \|\hat{\mathbf{s}}_t\| < +\infty$  almost surely and so the expectations  $\mathbb{E}_{t-1}[\mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \tilde{\mathbf{s}}_t]$  and  $\mathbb{E}_{t-1}[\mathbb{I}_{\{\hat{\gamma}_t \geq 1\}} \hat{\mathbf{s}}_t]$  are well defined. We also note that  $\gamma_t = \hat{\gamma}_t = \tilde{\gamma}_t$ , and  $\gamma_t, \mathbf{w}_t$ , and  $\mathbf{x}_t$  are all deterministic functions of the past ( $\mathbf{w}_t$  is the output of FTRL-prox, which is a deterministic function of the past).

By Lemma 7, there exists a random variable  $\Delta_t \geq 0$  satisfying  $\mathbb{E}_{t-1}[\Delta_t]$  and such that for all  $t \geq 1$ ,

$$\begin{aligned} \forall \mathbf{x} \in \mathcal{C}, \quad \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle &\leq \langle \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \tilde{\mathbf{v}}_t, \mathbf{w}_t - \mathbf{x} \rangle + (2\delta_t + \Delta_t) R \|\mathbf{g}_t\|, \\ &= \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \langle \mathbf{v}_t - \tilde{\mathbf{v}}_t, \mathbf{w}_t - \mathbf{x} \rangle \\ &\quad + (2\delta_t + \Delta_t) R \|\mathbf{g}_t\|. \end{aligned} \quad (91)$$

Fix  $\mathbf{x} \in \mathcal{C}$  and let  $X_t := \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \cdot \langle \mathbf{v}_t - \tilde{\mathbf{v}}_t, \mathbf{w}_t - \mathbf{x} \rangle$ . We start by bounding  $|X_t|$  from above. By the fact that  $\mathbf{x}_t \in \mathcal{C}$  and the  $B$ -Lipschitzness of  $\ell_t$ , we have

$$|X_t| \leq 2 \|\mathbf{g}_t\| R \|\mathbf{v}_t - \tilde{\mathbf{v}}_t\| \leq 2RB \cdot (\|\mathbf{v}_t\| + \|\tilde{\mathbf{v}}_t\|) \leq 4dRB \cdot (\kappa + \delta), \quad (92)$$

where the last inequality follows by the fact that  $\mathbf{v}_t$  [resp.  $\tilde{\mathbf{v}}_t$ ] has the same distribution as  $\mathbb{I}_{\{\hat{\gamma}_t \geq 1\}} \hat{\mathbf{s}}_t$  [resp.  $\mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \tilde{\mathbf{s}}_t$ ], and  $\|\hat{\mathbf{s}}_t \cdot \mathbb{I}_{\{\hat{\gamma}_t \geq 1\}}\| \leq d \cdot (\delta/R + 1/r)$  by Lemma 32 [resp.  $\|\tilde{\mathbf{s}}_t \cdot \mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}}\| \leq$

$d \cdot (\delta/R + 1/r)$  by Proposition 11]. We now show that  $\mathbb{E}_{t-1}[X_t] = 0$ . Using the definition of  $X_t$ , we have

$$\begin{aligned}\mathbb{E}_{t-1}[X_t] &= \mathbb{E}_{t-1}[\mathbb{I}_{\{\tilde{\gamma}_t < 1\}} \cdot \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \cdot \langle \mathbf{v}_t - \tilde{\mathbf{v}}_t, \mathbf{w}_t - \mathbf{x} \rangle] \\ &\quad + \mathbb{E}_{t-1}[\mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \cdot \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \cdot \langle \mathbf{v}_t - \tilde{\mathbf{v}}_t, \mathbf{w}_t - \mathbf{x} \rangle], \\ &= \mathbb{E}_{t-1}[\mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \cdot \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \cdot \langle \hat{\mathbf{s}}_t - \tilde{\mathbf{s}}_t, \mathbf{w}_t - \mathbf{x} \rangle] = 0,\end{aligned}\tag{93}$$

where the equalities in (93) follow by the facts that  $\tilde{\gamma}_t = \hat{\gamma}_t = \gamma_t$  (by Lemma 32);  $\mathbf{v}_t = \hat{\mathbf{s}}_t = \mathbf{0}$  if  $\gamma_t \geq 1$ ; and that  $\mathbb{E}_{t-1}[\mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \cdot (\hat{\mathbf{s}}_t - \tilde{\mathbf{s}}_t)] = 0$  (by Lemma 32). Using that  $\mathbf{x}_t$  is a deterministic function of the past, we get

$$\begin{aligned}\mathbb{E}_{t-1}[X_t^2] &\leq \|\mathbf{x}_t\|^2 \|\mathbf{g}_t\|^2 \mathbb{E}_{t-1}[\|\mathbf{v}_t - \tilde{\mathbf{v}}_t\|^2 \cdot \|\mathbf{w}_t - \mathbf{x}\|^2], \\ &\leq 8R^4 \|\mathbf{g}_t\|^2 \mathbb{E}_{t-1}[(\|\mathbf{v}_t\|^2 + \|\tilde{\mathbf{v}}_t\|^2)], \quad (\mathbf{x}_t \in \mathcal{C}) \\ &\leq 8R^4 \|\mathbf{g}_t\|^2 \mathbb{E}_{t-1}[\mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \cdot (\|\hat{\mathbf{s}}_t\|^2 + \|\tilde{\mathbf{s}}_t\|^2)], \\ &\leq 8R^2 B^2 d(\kappa + \delta)^2.\end{aligned}\tag{94}$$

where the last inequality follows by Lemma 32 and Proposition 11. Thus, by (91), (92), (93), (94), and Freedman's inequality (Theorem 40), we have, for any  $\rho \in (0, 1)$ , with probability at least  $1 - \rho$ ,

$$\begin{aligned}&\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \\ &= \sum_{t=1}^T \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + \sum_{t=1}^T (2\delta_t + \Delta_t) R \|\mathbf{g}_t\| + 4RB(\kappa + \delta) \sqrt{2dT \ln \rho^{-1}} + 4dRB \cdot (\kappa + \delta), \\ &\leq 2R \sqrt{2 \sum_{t=1}^T \|\tilde{\mathbf{g}}_t\|^2} + \sum_{t=1}^T (2\delta_t + \Delta_t) R \|\mathbf{g}_t\| + 4RB(\kappa + \delta) \sqrt{2dT \ln \rho^{-1}} + 4dRB \cdot (\kappa + \delta),\end{aligned}\tag{95}$$

where the last inequality follows by the regret bound FTRL-prox in Proposition 25. We now bound the first two terms on the RHS of (95). We start with  $\sum_{t=1}^T \|\tilde{\mathbf{g}}_t\|^2$ . With  $Y_t := \|\tilde{\mathbf{g}}_t\|^2$ , we have

$$\begin{aligned}\mathbb{E}_{t-1}[Y_t] &\leq \mathbb{E}_{t-1}[2\|\mathbf{g}_t\|^2 + 2\|\mathbf{g}_t\|^2 \|\mathbf{x}_t\|^2 \|\mathbf{v}_t\|^2], \\ &\leq 2B^2 R^2 \mathbb{E}_{t-1}[\mathbb{I}_{\{\tilde{\gamma}_t < 1\}} \cdot \|\mathbf{v}_t\|^2] + 2B^2 R^2 \mathbb{E}_{t-1}[\mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \cdot \|\hat{\mathbf{s}}_t\|^2] + 2\|\mathbf{g}_t\|^2, \quad (\mathbf{x}_t \in \mathcal{C}) \\ &\leq 2dB^2(\kappa + \delta)^2 + 2B^2,\end{aligned}\tag{96}$$

where the last inequality follows by Lemma 32. Similarly, we also have

$$\begin{aligned}\mathbb{E}_{t-1}[Y_t^2] &\leq \mathbb{E}_{t-1}[8\|\mathbf{g}_t\|^4 + 8\|\mathbf{g}_t\|^4 \|\mathbf{x}_t\|^4 \|\mathbf{v}_t\|^4], \\ &= 8B^4 R^4 \mathbb{E}_{t-1}[\|\mathbf{v}_t\|^4] + 8B^4, \quad (\mathbf{x}_t \in \mathcal{C}) \\ &= 8B^4 R^4 \mathbb{E}_{t-1}[\mathbb{I}_{\{\tilde{\gamma}_t < 1\}} \cdot \|\mathbf{v}_t\|^4] + 8B^4 R^4 \mathbb{E}_{t-1}[\mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \cdot \|\hat{\mathbf{s}}_t\|^4] + 8B^4, \\ &\leq 8d^3 B^4 (\kappa + \delta)^4 + 8B^4,\end{aligned}\tag{97}$$

where the last inequality follows by Lemma 32. Also, we have

$$|Y_t| \leq \|\mathbf{g}_t\| + \|\mathbf{g}_t\| \|\mathbf{x}_t\| \|\mathbf{v}_t\| \leq B \cdot (1 + R\|\mathbf{v}_t\|) \leq B \cdot (1 + d \cdot (\kappa + \delta)),\tag{98}$$

where the last inequality follows by the fact that  $\mathbf{v}_t$  has the same distribution as  $\mathbb{I}_{\{\hat{\gamma}_t \geq 1\}} \hat{\mathbf{s}}_t$ , and  $\|\hat{\mathbf{s}}_t\| \leq d \cdot (\delta/R + 1/r)$  by Lemma 32. By combining (96), (97), and (98), and applying Freedman's inequality (Theorem 40), we get for all  $\rho > 0$ , with probability at least  $1 - \rho$ ,

$$\begin{aligned} \sum_{t=1}^T \|\tilde{\mathbf{g}}_t\|^2 &\leq \sum_{t=1}^T \mathbb{E}_{t-1}[\|\tilde{\mathbf{g}}_t\|^2] + 4\sqrt{2B^4(1 + d^3(\kappa + \delta)^4)T \ln(1/\rho)} + B^2(1 + d(\kappa + \delta))^2 \ln(1/\delta), \\ &\leq 4B^2d(\kappa + \delta)^2T + 8d^{3/2}B^2(\kappa + \delta)^2\sqrt{T \ln(1/\rho)} + B^2(1 + d(\kappa + \delta))^2 \ln(1/\delta), \\ &\leq 12B^2d(\kappa + \delta)^2T + 4B^2d^2(\kappa + \delta)^2 \ln(1/\delta), \end{aligned} \quad (99)$$

where the last inequality follows by the fact that  $d \ln(1/\rho) \leq T$ . Finally, by Markov's inequality and the fact that  $\ell_t$  is  $B$ -Lipschitz, we have, for all  $\rho$ ,

$$\mathbb{P}\left[\sum_{t=1}^T \Delta_t \|\mathbf{g}_t\| \geq 2B\delta\sqrt{T}\rho^{-1}\right] \leq \frac{\rho \sum_{t=1}^T \mathbb{E}[\Delta_t \|\mathbf{g}_t\|]}{2\delta B\sqrt{T}} \leq \frac{\rho \sum_{t=1}^T \mathbb{E}[\Delta_t]}{2\delta\sqrt{T}} \leq \frac{\rho \sum_{t=1}^T \delta_t}{2\delta} \leq \rho, \quad (100)$$

where the last inequality follows the fact that  $\delta_t = \delta/t^{-1/2}$  and  $\sum_{t=1}^T 1/t^{-1/2} \leq 2\sqrt{T}$ , for all  $T \geq 1$ . By combining (95), (99), and (100), and a union bound, we get, with probability at least  $1 - 3\rho$ ,

$$\begin{aligned} &\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \\ &\leq 4RB(\kappa + \delta)\sqrt{6dT + 2d \ln \rho^{-1}} + 6RB\delta\sqrt{T}/\rho + 4RB(\kappa + \delta)\sqrt{2dT \ln \rho^{-1}} + 4dRB \cdot (\kappa + \delta), \\ &\leq 8RB(\kappa + \delta)\sqrt{dT(3 + \ln \rho^{-1}) + d \ln \rho^{-1}} + 2dBR \cdot (2\kappa + \delta \cdot (2 + 3\sqrt{T}/\rho)), \end{aligned}$$

where the last inequality follows by the fact that  $\sqrt{a} + \sqrt{b} \leq \sqrt{2a + 2b}$ . This completes the proof. ■

### I.5. Proof of Theorem 17 (The Strongly Convex Case)

Before proving Theorem 17, we present a result that may be viewed as an extension of Lemma 7:

**Lemma 33** *Let  $\mathbf{w}_t, \tilde{\mathbf{g}}_t$ , and  $\mathbf{x}_t$  be as in Algorithm 1 with  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{\mathcal{C}^\circ}$ , any tolerance sequence  $(\delta_s) \subset (0, 1/3)$ , and any OCO subroutine  $A$  defined on  $\mathcal{B}(R)$  ( $R$  as in (1)). Then, for all  $t \geq 1$  and all  $\mathbf{x} \in \mathcal{C}$ , we have*

$$\langle \tilde{\mathbf{g}}_t, \mathbf{x}_t \rangle \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t \rangle + (2\delta_t + \Delta_t)R\|\mathbf{g}_t\|,$$

where  $\Delta_t \in [0, 15^2d^4\kappa^3\delta_t^{-2}]$  is the random variable in Lemma 7, which satisfies  $\mathbb{E}_{t-1}[\Delta_t] \leq \delta_t$ .

**Proof** Instantiating (5) in Lemma 7 with  $\mathbf{x} = \mathbf{x}_t$ , we get

$$\langle \tilde{\mathbf{g}}_t, \mathbf{w}_t \rangle \geq \tilde{\ell}_t(\mathbf{w}_t) - \tilde{\ell}_t(\mathbf{x}_t) + \langle \tilde{\mathbf{g}}_t, \mathbf{x}_t \rangle - (\delta_t + \Delta_t)R\|\mathbf{g}_t\|. \quad (101)$$

Now, to get the desired result, we need to show that  $\tilde{\ell}_t(\mathbf{w}_t) - \tilde{\ell}_t(\mathbf{x}_t) \geq -\delta_t R\|\mathbf{g}_t\|$ . When  $\gamma_t < 1$ , we have  $\mathbf{x}_t = \mathbf{w}_t$  and so

$$\tilde{\ell}_t(\mathbf{x}_t) = \tilde{\ell}_t(\mathbf{w}_t). \quad [\text{case where } \gamma_t < 1] \quad (102)$$

Now, suppose that  $\gamma_t \geq 1$  and  $\langle \mathbf{g}_t, \mathbf{w}_t \rangle \geq 0$ . By the fact that  $\mathbf{x}_t = \mathbf{w}_t/\gamma_t$ , we get

$$\tilde{\ell}_t(\mathbf{x}_t) = \langle \mathbf{g}_t, \mathbf{x}_t \rangle \leq \langle \mathbf{g}_t, \mathbf{w}_t \rangle = \tilde{\ell}_t(\mathbf{w}_t). \quad [\text{case where } \gamma_t \geq 1, \langle \mathbf{g}_t, \mathbf{w}_t \rangle \geq 0] \quad (103)$$

Now suppose that  $\gamma_t \geq 1$ ,  $\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0$ , and  $\mathbf{w}_t \in \mathcal{C}$ . We note that this implies that  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \leq 1$  and so  $S_{\mathcal{C}}(\mathbf{w}_t) = 0$  (Lemma 6). Thus,  $\tilde{\ell}_t(\mathbf{w}_t) = \langle \mathbf{g}_t, \mathbf{w}_t \rangle$ . On the other hand, by Lemma 10 we have  $\gamma_t \leq \gamma_{\mathcal{C}}(\mathbf{w}_t) + \delta_t \leq 1 + \delta_t$ , and so since  $\mathbf{x}_t = \mathbf{w}_t/\gamma_t$ , we have

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle \leq \langle \mathbf{g}_t, \mathbf{w}_t \rangle / (1 + \delta_t) = \langle \mathbf{g}_t, \mathbf{w}_t \rangle - \delta_t \langle \mathbf{g}_t, \mathbf{w}_t \rangle / (1 + \delta_t).$$

Thus, since  $\tilde{\ell}_t(\mathbf{w}_t) = \langle \mathbf{g}_t, \mathbf{w}_t \rangle$  and  $\tilde{\ell}_t(\mathbf{x}_t) = \langle \mathbf{g}_t, \mathbf{x}_t \rangle$  (because  $\mathbf{x}_t \in \mathcal{C}$ ), the previous display implies that

$$\tilde{\ell}_t(\mathbf{x}_t) = \langle \mathbf{g}_t, \mathbf{x}_t \rangle \leq \tilde{\ell}_t(\mathbf{w}_t) + \delta_t R \|\mathbf{g}_t\|. \quad [\text{case where } \gamma \geq 1, \langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0, \mathbf{w}_t \in \mathcal{C}] \quad (104)$$

We now consider the last case where  $\gamma_t \geq 1$ ,  $\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0$ , and  $\mathbf{w}_t \notin \mathcal{C}$ . We note that this implies that  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \geq 1$ . By the fact that  $\gamma_t \leq \gamma_{\mathcal{C}}(\mathbf{w}_t) + \delta_t$  (Lemma 10), we have

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle \leq \frac{\langle \mathbf{g}_t, \mathbf{w}_t \rangle}{\gamma_{\mathcal{C}}(\mathbf{w}_t) + \delta_t} = \frac{\langle \mathbf{g}_t, \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t) \rangle}{1 + \delta_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)} = \langle \mathbf{g}_t, \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t) \rangle - \frac{\delta_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)}{1 + \delta_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)} \langle \mathbf{g}_t, \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t) \rangle, \quad (105)$$

Thus, since  $S_{\mathcal{C}}(\mathbf{w}_t) = \gamma_{\mathcal{C}}(\mathbf{w}_t) - 1$  (by Lemma 6 and  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \geq 1$ ), we get

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle \cdot S_{\mathcal{C}}(\mathbf{w}_t) \leq \langle \mathbf{g}_t, \mathbf{w}_t \rangle - \frac{\delta_t \langle \mathbf{g}_t, \mathbf{w}_t \rangle}{1 + \delta_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)} - \langle \mathbf{g}_t, \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t) \rangle + \frac{\delta_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)}{1 + \delta_t / \gamma_{\mathcal{C}}(\mathbf{w}_t)} \langle \mathbf{g}_t, \mathbf{w}_t / \gamma_{\mathcal{C}}(\mathbf{w}_t) \rangle.$$

Adding this together with (105) and using the fact that  $\|\mathbf{w}_t\| \leq R$ , we get

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle + \langle \mathbf{g}_t, \mathbf{x}_t \rangle S_{\mathcal{C}}(\mathbf{w}_t) \leq \langle \mathbf{g}_t, \mathbf{w}_t \rangle + \delta_t R \|\mathbf{g}_t\|,$$

and so after rearranging and using that  $\tilde{\ell}_t(\mathbf{x}_t) = \langle \mathbf{x}_t, \mathbf{g}_t \rangle$  (since  $\mathbf{x}_t \in \mathcal{C}$ ) and  $\tilde{\ell}_t(\mathbf{w}_t) = \langle \mathbf{g}_t, \mathbf{w}_t \rangle - \langle \mathbf{g}_t, \mathbf{x}_t \rangle S_{\mathcal{C}}(\mathbf{w}_t)$ , we get

$$\langle \mathbf{g}_t, \mathbf{x}_t \rangle \leq \tilde{\ell}_t(\mathbf{w}_t) + \delta_t R \|\mathbf{g}_t\|. \quad [\text{case where } \gamma \geq 1, \langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0, \mathbf{w}_t \notin \mathcal{C}] \quad (106)$$

By combining, (102), (103), (104), and (106), we obtain:

$$\tilde{\ell}_t(\mathbf{x}_t) \leq \tilde{\ell}_t(\mathbf{w}_t) + \delta_t R \|\mathbf{g}_t\|.$$

Combining this with (101), we get

$$\langle \tilde{\mathbf{g}}_t, \mathbf{x}_t \rangle \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t \rangle + (2\delta_t + \Delta_t) R \|\mathbf{g}_t\|.$$

■

We now present the versions of Lemmas 7 and 33 when  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{1d, \mathcal{C}^\circ}$ :



**Lemma 34** Let  $\mathbf{w}_t, \tilde{\mathbf{g}}_t$ , and  $\mathbf{x}_t$  be as in Algorithm 1 with  $\mathcal{O}_{\mathcal{C}^\circ} \equiv \text{OPT}_{1d, \mathcal{C}^\circ}$ ; any tolerance sequence  $(\delta_s) \subset (0, 1/3)$ ; any tolerance sequence  $(\delta_t) \subset (0, 1/3)$ ; and any OCO subroutine  $A$  defined on  $\mathcal{B}(R)$  ( $R$  as in (1)). Then, for all  $t \geq 1$ ,  $\mathbf{x}_t \in \mathcal{C}$  and there exists a family of random variables  $\{\xi_t(\mathbf{x}) : \mathbf{x} \in \mathcal{C}\}$  such that for all  $\mathbf{x} \in \mathcal{C}$ ,  $\mathbb{E}_{t-1}[\xi_t(\mathbf{x})] = 0$ ,  $|\xi_t(\mathbf{x})| < 4d(\kappa + \delta_t)R\|\mathbf{g}_t\|$  almost surely, and

$$\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + (2\delta_t + \Delta_t)R\|\mathbf{g}_t\| + \xi_t(\mathbf{x}); \quad (107)$$

$$\langle \tilde{\mathbf{g}}_t, \mathbf{x}_t \rangle \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t \rangle + (2\delta_t + \Delta_t)R\|\mathbf{g}_t\| + \xi_t(\mathbf{x}_t), \quad (108)$$

where  $\Delta_t \in [0, 15^2 d^4 \kappa^3 \delta_t^{-2}]$  is a random variable satisfying  $\mathbb{E}_{t-1}[\Delta_t] \leq \delta_t$ . Furthermore,  $\|\tilde{\mathbf{g}}_t\| \leq (1 + d\kappa + d\Delta_t)\|\mathbf{g}_t\|$  and  $\mathbb{E}_{t-1}[\|\tilde{\mathbf{g}}_t\|^2] \leq 2(1 + d(\kappa + \delta_t)^2)\|\mathbf{g}_t\|^2$ .

**Proof** Let  $t \geq 1$ ;  $\mathbf{x} \in \mathcal{C}$ ;  $(\tilde{\gamma}_t, \tilde{\mathbf{s}}_t) = \text{OPT}_{\mathcal{C}^\circ}(\mathbf{w}_t; \delta_t)$ ; and  $(\hat{\gamma}_t, \hat{\mathbf{s}}_t) = \text{OPT}_{1d, \mathcal{C}^\circ}(\mathbf{w}_t; \delta_t)$ . Further, let  $\tilde{\mathbf{v}}_t := \mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \tilde{\mathbf{s}}_t$  and  $\xi_t(\mathbf{x}) := \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \langle \mathbf{v}_t - \tilde{\mathbf{v}}_t, \mathbf{w}_t - \mathbf{x} \rangle$ . By Lemma 7, there exists a random variable  $\Delta_t \geq 0$  satisfying  $\mathbb{E}_{t-1}[\Delta_t] \leq \delta_t$  and such that for all  $t \geq 1$  and  $\mathbf{x} \in \mathcal{C}$ ,

$$\begin{aligned} \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle &\leq \langle \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \tilde{\mathbf{v}}_t, \mathbf{w}_t - \mathbf{x} \rangle + (2\delta_t + \Delta_t)R\|\mathbf{g}_t\|, \\ &\leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + (2\delta_t + \Delta_t)R\|\mathbf{g}_t\| + \xi_t(\mathbf{x}). \end{aligned}$$

By Proposition 11 and Lemma 32, we have  $\|\tilde{\mathbf{s}}_t\| \vee \|\hat{\mathbf{s}}_t\| < d \cdot (1/R + \delta/r)$  whenever  $\hat{\gamma}_t = \tilde{\gamma}_t \geq 1$ . This implies that  $\|\tilde{\mathbf{v}}_t\| \vee \|\hat{\mathbf{v}}_t\| < d \cdot (1/R + \delta/r)$ , and so  $|\xi_t(\mathbf{x})| < 4d(\kappa + \delta)R\|\mathbf{g}_t\|$ . We now show that the conditional expectation of  $\xi_t(\mathbf{x})$  is zero. By Lemma 13, and the law of total expectation, we have

$$\begin{aligned} \mathbb{E}_{t-1}[\xi_t(\mathbf{x})] &\leq \mathbb{E}_{t-1}[\mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \langle \mathbf{v}_t - \tilde{\mathbf{v}}_t, \mathbf{w}_t - \mathbf{x} \rangle], \\ &= \mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \cdot \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \langle \mathbb{E}_{t-1}[\mathbf{v}_t - \tilde{\mathbf{v}}_t], \mathbf{w}_t - \mathbf{x} \rangle, \quad (109) \\ &= \mathbb{I}_{\{\tilde{\gamma}_t \geq 1\}} \cdot \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \langle \mathbb{E}_{t-1}[\mathbf{s}_t - \tilde{\mathbf{s}}_t], \mathbf{w}_t - \mathbf{x} \rangle, \\ &= 0. \quad (110) \end{aligned}$$

where (109) follows by the fact that  $\mathbf{w}_t, \gamma_t, \tilde{\gamma}_t$ , and  $\mathbf{g}_t$  are deterministic function of the past, and (110) follows by the fact  $\tilde{\gamma}_t = \hat{\gamma}_t$ ; and  $\mathbb{E}_{t-1}[\hat{\mathbf{s}}_t] = \mathbb{E}_{t-1}[\tilde{\mathbf{s}}_t]$ , when  $\tilde{\gamma}_t \geq 1$  (Lemma 13). This shows (107). We follow similar steps to show (108), but we use the result of Lemma 33 instead of Lemma 7. By Lemma 33, we have

$$\langle \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \tilde{\mathbf{v}}_t, \mathbf{x}_t \rangle \leq \langle \mathbf{g}_t - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \tilde{\mathbf{v}}_t, \mathbf{w}_t \rangle + (2\delta_t + \Delta_t)R\|\mathbf{g}_t\|,$$

and so by rearranging, we get

$$\langle \tilde{\mathbf{g}}_t, \mathbf{x}_t \rangle \leq \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t \rangle + (2\delta_t + \Delta_t)R\|\mathbf{g}_t\| + \xi_t(\mathbf{x}_t).$$

This shows (108). We now bound  $\|\tilde{\mathbf{g}}_t\|$ . When  $\tilde{\gamma}_t < 1$ , we have  $\mathbf{g}_t = \tilde{\mathbf{g}}_t$  and so  $\|\tilde{\mathbf{g}}_t\| \leq (1 + d\kappa + d\Delta_t)\|\mathbf{g}_t\|$  holds trivially. Now suppose that  $\tilde{\gamma}_t \geq 1$ . In this case,  $\mathbf{v}_t$  has the same distribution as  $\mathbb{I}_{\{\hat{\gamma}_t \geq 1\}} \hat{\mathbf{s}}_t$ , and so

$$\begin{aligned} \|\tilde{\mathbf{g}}_t\| &= \|\langle \mathbf{g}_t, \mathbf{w} \rangle - \mathbb{I}_{\langle \mathbf{g}_t, \mathbf{w}_t \rangle < 0} \langle \mathbf{g}_t, \mathbf{x}_t \rangle \mathbf{v}_t\| = \|\mathbf{g}_t\| + \|\mathbf{g}_t\| \frac{\|\mathbf{w}_t\|}{\gamma_t} \|\mathbf{v}_t\| \leq \|\mathbf{g}_t\| (1 + R\|\mathbf{v}_t\|) \\ &\leq (1 + d\Delta_t + d\kappa)\|\mathbf{g}_t\|, \end{aligned}$$

where the last inequality follows from the fact that  $\mathbb{I}_{\{\hat{\gamma}_t \geq 1\}} \|\hat{\mathbf{s}}_t\| \leq d \cdot (\Delta_t/R + 1/r)$ , by Lemma 32. From the penultimate inequality in the above display and Lemma 13, we also have

$$\mathbb{E}_{t-1}[\|\tilde{\mathbf{g}}_t\|^2] \leq 2\|\mathbf{g}_t\|^2(1 + R^2\mathbb{E}_{t-1}[\|\hat{\mathbf{s}}_t\|^2]) \leq 2\|\mathbf{g}_t\|^2(1 + d \cdot (\kappa + \delta_t)^2).$$

Finally, the fact that  $\mathbf{x}_t \in \mathcal{C}$  follows from Theorem 15. In fact, looking at the proof of Theorem 15 it is clear that the fact that  $(\mathbf{x}_t) \subset \mathcal{C}$  only uses the fact that the subroutine A outputs iterates in  $(\mathbf{x}_t)$  in  $\mathcal{B}(R)$ .  $\blacksquare$

We now present the proof of Theorem 17. We note that the proof follows from that of (Cutkosky and Orabona, 2018, Theorem 7) with small modifications to account for the differences between their constrained-to-unconstrained reduction and our projection-free reduction (Algorithm 1), as well as the differences described in Section C.2 to achieve scale-invariance.

**Proof of Theorem 17.** First of all, note that  $(\mathbf{x}_t) \subset \mathcal{C}$  follows directly from Lemma 34. Next, we will instantiate Theorem 35 below with  $\delta_t = \delta/t^2$  for all  $t \geq 1$  and some  $\delta \in (0, 1/3)$ . By Lemma 34 and the fact that the losses are  $B$ -Lipschitz, the variables  $\zeta_T, Z_T, M_T$ , and  $N_T$  in Theorem 35 below satisfy

- $\mathbb{E}[\zeta_T] \leq 1 + d\kappa + 2d\delta_T$ , where we used the fact that  $\mathbb{E}_{t-1}[\Delta_t] \leq \delta_t$  and  $\sum_{t=1}^{\infty} 1/t^2 \leq 2$ .
- $\zeta_T \leq 16^2 d^5 \kappa^3 / \delta_T^2$ , where we used the fact that  $\Delta_t \leq 15^2 d^4 \kappa^3 / \delta_t^2, \forall t$ .
- $Z_T \leq \epsilon^2 + B^2(16^4 d^{10} \kappa^6 / \delta_T^4)(1 + T) \leq B^2(16^4 d^{10} \kappa^6 / \delta_T^4)(2 + T)$ .
- $M_T \leq 2\sqrt{\ln_+ \left( 2(16^6 d^{15} \kappa^9 / \delta_T^6)(2T + 1)^2 RB^2 / \epsilon^2 \right)}$ .
- $\mathbb{E}[N_T] \leq \mathbb{E}[4(1 + d\kappa + 2d\delta)B \ln_+ \left( 4(16^6 d^{15} \kappa^9 / \delta_T^6)(2T + 1)^2 RB^2 / \epsilon^2 \right) + \epsilon/R + (1 + 2d\kappa + 8d\delta)B]$ . This follows by the facts that  $\mathbb{E}[\zeta_T] \leq 1 + d\kappa + 2d\delta$  and  $\zeta_T \leq 16^2 d^5 \kappa^3 / \delta_T^2$ .
- $\mathbb{E}_{t-1}[\|\tilde{\mathbf{g}}_t\|^2] \leq 2(1 + d(\kappa + \delta_t)^2)\|\mathbf{g}_t\|^2$  by Lemma 34.

Thus, for  $\nu := 1/(R \wedge 1) + \kappa + 2\delta$ , there exists  $W_T = O(\ln(e + \kappa dRTB/(\delta\epsilon)))$  such that

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \right] &\leq W_T \sqrt{\mathbb{E} \left[ \sum_{t=1}^T \|\tilde{\mathbf{g}}_t\|^2 \|\mathbf{x}_t - \mathbf{x}\|^2 \right]} + d\nu RBW_T^2, \\ &= W_T \sqrt{\mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}_{t-1}[\|\tilde{\mathbf{g}}_t\|^2 \|\mathbf{x}_t - \mathbf{x}\|^2] \right]} + d\nu RBW_T^2, \\ &\leq W_T \sqrt{\mathbb{E} \left[ 2d\nu^2 \sum_{t=1}^T \|\mathbf{g}_t\|^2 \|\mathbf{x}_t - \mathbf{x}\|^2 \right]} + d\nu RBW_T^2, \\ &\leq U_T \sqrt{\mathbb{E} \left[ \sum_{t=1}^T \|\mathbf{g}_t\|^2 \|\mathbf{x}_t - \mathbf{x}\|^2 \right]} + RBW_T^2/\nu, \end{aligned} \tag{111}$$

where  $U_T = O(\nu d^{1/2} \ln(e + \kappa d R T B / (\delta \epsilon)))$ . Thus, for general convex functions, we have, in expectation

$$\sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{x})) \leq U_T R B \sqrt{T} + R B U_T^2 / \nu.$$

For  $\mu$ -strongly convex functions ( $\ell_t$ ), we have

$$\begin{aligned} & \mathbb{E} \left[ \sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{x})) \right] \\ & \leq \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \right] - \mathbb{E} \left[ \frac{\mu}{2} \sum_{t=1}^T \|\mathbf{x}_t - \mathbf{x}\|^2 \right], \\ & \leq U_T \sqrt{\mathbb{E} \left[ \sum_{t=1}^T \|\mathbf{g}_t\|^2 \|\mathbf{x}_t - \mathbf{x}\|^2 \right]} + R B U_T^2 / \nu - \mathbb{E} \left[ \frac{\mu}{2} \sum_{t=1}^T \|\mathbf{x}_t - \mathbf{x}\|^2 \right], \end{aligned} \quad (112)$$

$$\begin{aligned} & \leq \inf_{\eta > 0} \left\{ \mathbb{E} \left[ \sum_{t=1}^T \|\mathbf{g}_t\|^2 \|\mathbf{x}_t - \mathbf{x}\|^2 / \eta \right] + \eta U_T^2 / 4 \right\} - \mathbb{E} \left[ \frac{\mu}{2} \sum_{t=1}^T \|\mathbf{x}_t - \mathbf{x}\|^2 \right] + \nu R B U_T^2 / \nu, \\ & \leq \frac{\mu}{2} \mathbb{E} \left[ \sum_{t=1}^T \frac{\|\mathbf{g}_t\|^2}{B^2} \|\mathbf{x}_t - \mathbf{x}\|^2 \right] + \frac{B^2 U_T^2}{2\mu} - \mathbb{E} \left[ \frac{\mu}{2} \sum_{t=1}^T \|\mathbf{x}_t - \mathbf{x}\|^2 \right] + R B U_T^2 / \nu, \\ & \leq B^2 U_T^2 / (2\mu) + R B U_T^2 / \nu, \end{aligned} \quad (113)$$

where (112) follows by (111) and (113) follows by setting  $\eta = 2B^2/\mu$ .  $\blacksquare$

**Theorem 35** *Let  $\delta \in (0, 1/3)$  and  $\kappa := R/r$ , with  $r$  and  $R$  as in (1). Suppose that Algorithm 1 is run with  $\mathcal{O}_{C^\circ} \equiv \text{OPT}_{\text{Id}, C^\circ}$ ;  $\delta_t = \delta/t^2$ ,  $\forall t \geq 1$ ; and sub-routine A set to Alg. 5 with parameter  $\epsilon > 0$ . Then, for any adversarial sequence of convex losses ( $\ell_t$ ) on  $\mathcal{C}$ , the iterates ( $\mathbf{x}_t$ ) of Alg. 1 in response to ( $\ell_t$ ) satisfy,*

$$\begin{aligned} \forall T \geq 1, \forall \mathbf{x} \in \mathcal{C}, \quad \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle & \leq M_T \sqrt{2 \left[ (\epsilon^2 + \zeta_T B_T^2) \|\mathbf{x}\|^2 + \sum_{t=1}^T \|\tilde{\mathbf{g}}_t\|^2 \|\mathbf{x}_t - \mathbf{x}\|^2 \right]} \cdot \ln \frac{Z_T}{\epsilon^2} \\ & \quad + 2R N_T \cdot \left( 1 + \ln \frac{Z_T}{\epsilon^2} \right) + \Xi_T(\mathbf{x}), \end{aligned}$$

where  $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$ ,  $\forall t$ ,  $B_T := \epsilon \vee \max_{t \in [T]} \|\mathbf{g}_t\|$ , and  $\zeta_T, Z_T, M_T, \Xi_T(\mathbf{x})$  and  $N_T$  are such that:

- $\zeta_T := 1 + d\kappa + d \max_{t \in [T]} \Delta_t$ , with  $\Delta_t \geq 0$  a non-negative random variable satisfying  $\mathbb{E}_{t-1}[\Delta_t] \leq \delta_t$ .
- $Z_T \leq \epsilon^2 + \tilde{B}_T^2 + \sum_{t=1}^T \|\tilde{\mathbf{g}}_t\|^2 \leq \epsilon^2 + \zeta_T^2 B_T^2 + \zeta_T^2 \sum_{t=1}^T \|\mathbf{g}_t\|^2$ .
- $M_T := 2\sqrt{\ln_+ (2\zeta_t^2 (R + 2\zeta_t R T) Q_T / \epsilon^2)}$ , and  $Q_t := \epsilon^2 + \sum_{i=1}^t \|\mathbf{g}_i\|^2$ .

- $\Xi_T(\mathbf{x}) := \sum_{t=1}^T \xi_t(\mathbf{x})$ , where  $\{\xi_t(\mathbf{x}) : \mathbf{x} \in \mathcal{C}\}$  is a family of random variables satisfying for all  $\mathbf{x} \in \mathcal{C}$ ,  $\mathbb{E}_{t-1}[\xi_t(\mathbf{x})] = 0$ .
- $N_T := 4\zeta_t B_T \ln_+ (4\zeta_t^2 B_T (R + 2\zeta_t R T) \sqrt{Q_T}/\epsilon^2) + B_T + \frac{\epsilon}{R} + \sum_{t=1}^T (2\delta_t + \Delta_t) \|\mathbf{g}_t\| + 2d(\kappa + \delta) B_T$ .

**Proof** Define  $\zeta_t := 1 + d\kappa + d \max_{s \leq t} \Delta_s$ , where  $\kappa = R/r$  and  $\Delta_t$  is the same random variable as in Lemma 34. Note that  $\zeta_t \geq 0$  and  $\mathbb{E}[\zeta_t] \leq 1 + d\kappa + d \sum_{s=1}^t \delta_s$ . For any  $t \geq 1$ , consider the random vector  $\mathbf{X}_t$  that takes value  $\mathbf{x}_i$  for  $s \leq t$  with probability proportional to  $\|\hat{\mathbf{g}}_i\|^2$ , and value  $\mathbf{0}$  with probability proportional to  $\epsilon^2 + \tilde{B}_t^2$ , where  $\tilde{B}_t := \epsilon \vee \max_{i \in [t]} \|\hat{\mathbf{g}}_i\|$  and  $\hat{\mathbf{g}}_t := \tilde{\mathbf{g}}_t \cdot \tilde{B}_{t-1}/\tilde{B}_t$  (see Algorithm 5). Moving forward, we define

$$z_t := \|\hat{\mathbf{g}}_t\|^2, \quad \text{and} \quad z_{0,t} := \epsilon^2 + \tilde{B}_t^2,$$

so that  $Z_t := z_{0,t} + \sum_{i=1}^t z_i = \epsilon^2 + \tilde{B}_t^2 + \sum_{i=1}^t \|\hat{\mathbf{g}}_i\|^2$ . We make the following definitions/observations:

- We define  $V_T(\mathbf{x}) = z_{0,T} \|\mathbf{x}\|^2 + \sum_{t=1}^T z_t \|\mathbf{x}_t - \mathbf{x}\|^2 = Z_T \cdot \mathbb{E}[\|\mathbf{X}_T - \mathbf{x}\|^2]$ .
- $Z_t \leq \zeta_t^2 B_T^2 \cdot (T + 1) + \epsilon^2$ , for all  $t \in [T]$ , which follows from Lemma 34.
- $\bar{\mathbf{x}}_t = \mathbb{E}[\mathbf{X}_t] = Z_t^{-1} \cdot \sum_{i=1}^t z_i \cdot \mathbf{x}_i = \mathbf{v}_t / Z_t$ , where  $(\mathbf{v}_i)$  are as in Algorithm 5.
- $\sigma_t^2 := (z_{0,t} \|\bar{\mathbf{x}}_t\|^2 + \sum_{i=1}^t z_i \cdot \|\mathbf{x}_i - \bar{\mathbf{x}}_t\|^2) / Z_t$  so that  $\sigma_t^2 = \mathbb{E}[\|\mathbf{X}_t - \bar{\mathbf{x}}_t\|^2]$ .

To prove the theorem, we are going to show for any  $\mathbf{x} \in \mathcal{C}$ ,

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq M_T \sqrt{Z_T \cdot \|\mathbf{x} - \mathbf{x}_T\|^2} + M_T \sqrt{\sigma_T^2 \cdot Z_T \cdot \ln \frac{Z_T}{\epsilon^2}} + 2RN_T \cdot \left(1 + \ln \frac{Z_T}{\epsilon^2}\right) + \Xi_T(\mathbf{x}), \quad (114)$$

where  $M_T$  and  $N_T$  are as in (116) and (117) below, respectively, which implies the desired bound by a bias-variance decomposition:

$$Z_T \cdot \|\mathbf{x} - \bar{\mathbf{x}}_T\|^2 + Z_T \cdot \sigma_T^2 = Z_T \cdot \mathbb{E}[\|\mathbf{X}_T - \mathbf{x}\|^2] = V_T(\mathbf{x}).$$

In particular, combining this with (114) and using the fact that  $\sqrt{x} + \sqrt{y} \leq \sqrt{2(x+y)}$  for all  $x, y > 0$ , we get

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq M_T \sqrt{2V_T(\mathbf{x}) \cdot \|\mathbf{x} - \mathbf{x}_T\|^2 \cdot \ln \frac{Z_T}{\epsilon^2}} + 2RN_T \cdot \left(1 + \ln \frac{Z_T}{\epsilon^2}\right) + \Xi_T(\mathbf{x}).$$

To get to (114), first observe that for all  $\mathbf{x} \in \mathcal{C}$ , the linearized regret of Algorithm 5 satisfies

$$\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle - \sum_{t=1}^T \langle \bar{\mathbf{g}}_t, \mathbf{x}_t - \mathbf{x} \rangle = \sum_{t=1}^T \langle \mathbf{g}_t - \bar{\mathbf{g}}_t, \mathbf{x}_t - \mathbf{x} \rangle \stackrel{(*)}{\leq} 2R \sum_{t=1}^T (B_t - B_{t-1}) \leq 2RB_T,$$

where (\*) following by Cauchy Schwarz. Using this together with (107) (multiplied by  $\tilde{B}_{t-1}/\tilde{B}_t$ ) and the clipped regret of FreeGrad in Proposition 26, we have, for all  $\mathbf{x} \in \mathcal{C}$ ,

$$\begin{aligned}
 & \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle - 2RB_T \\
 & \leq \sum_{t=1}^T \langle \hat{\mathbf{g}}_t, \mathbf{x}_t - \mathbf{x} \rangle, \\
 & \leq \sum_{t=1}^T \langle \hat{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle + \sum_{t=1}^T (2\delta_t + \Delta_t) R \|\mathbf{g}_t\| + \sum_{t=1}^T \xi_t(\mathbf{x}), \\
 & = \sum_{t=1}^T \langle \hat{\mathbf{g}}_t, \mathbf{u}_t - (\mathbf{x} - \bar{\mathbf{x}}_T) \rangle + \sum_{t=1}^T \langle \hat{\mathbf{g}}_t, \bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_T \rangle + \sum_{t=1}^T (\xi_t(\mathbf{x}) + (2\delta_t + \Delta_t) R \|\mathbf{g}_t\|), \\
 & = M_T \sqrt{Z_T \cdot \|\mathbf{x} - \bar{\mathbf{x}}_T\|^2} + \sum_{t=1}^T \langle \hat{\mathbf{g}}_t, \bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_T \rangle + 2RN_T + \Xi_T(\mathbf{x}), \tag{115}
 \end{aligned}$$

where  $\Xi_T(\mathbf{x}) := \sum_{t=1}^T \xi_t(\mathbf{x})$ . Note that the first term on the RHS of (115) is exactly what we want, so we only have to derive an upper bound on the second one. This is readily done through Lemma 36 that immediately gives us the stated result.  $\blacksquare$

**Lemma 36** *Under the hypotheses of Theorem 35, we have*

$$\begin{aligned}
 \sum_{t=1}^T \langle \hat{\mathbf{g}}_t, \bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_T \rangle & \leq M_T \sigma_T \sqrt{Z_T \ln \frac{Z_T}{\epsilon^2}} + 2RN_T \cdot \ln \frac{Z_T}{\epsilon^2}, \text{ where} \\
 M_T & := 2 \sqrt{\ln_+ \left( \frac{2\zeta_t^2 (R + 2\zeta_t RT) Q_T}{\epsilon^2} \right)}; \quad Q_T := \epsilon^2 + \sum_{t=1}^T \|\mathbf{g}_t\|^2; \quad \text{and} \tag{116}
 \end{aligned}$$

$$N_T := 4\zeta_t B_T \ln_+ \left( \frac{4\zeta_t^2 B_T (R + 2\zeta_t RT) \sqrt{Q_T}}{\epsilon^2} \right) + B_T + \frac{\epsilon}{R} + \sum_{t=1}^T (2\delta_t + \Delta_t) \|\mathbf{g}_t\| + 2d(\kappa + \delta) B_T. \tag{117}$$

**Proof** We have that

$$\sum_{i=1}^t \langle \hat{\mathbf{g}}_i, \bar{\mathbf{x}}_{i-1} - \bar{\mathbf{x}}_t \rangle - \sum_{i=1}^{t-1} \langle \hat{\mathbf{g}}_i, \bar{\mathbf{x}}_{i-1} - \bar{\mathbf{x}}_{t-1} \rangle = \left\langle \sum_{i=1}^t \hat{\mathbf{g}}_i, \bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_t \right\rangle.$$

The telescoping sum gives us

$$\sum_{t=1}^T \langle \hat{\mathbf{g}}_t, \bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_T \rangle = \sum_{t=1}^T \left\langle \sum_{i=1}^t \hat{\mathbf{g}}_i, \bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_t \right\rangle \leq \sum_{t=1}^T \left\| \sum_{i=1}^t \hat{\mathbf{g}}_i \right\| \|\bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_t\|.$$

So in order to bound  $\sum_{t=1}^T \langle \hat{\mathbf{g}}_t, \bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_T \rangle$ , it suffices to bound  $\left\| \sum_{i=1}^t \hat{\mathbf{g}}_i \right\| \|\bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_t\|$  by a sufficiently small value. First, we will tackle  $\left\| \sum_{i=1}^t \hat{\mathbf{g}}_i \right\|$ . By Lemma 34 (in particular (108)) and the

fact that  $\hat{\mathbf{g}}_t = \tilde{\mathbf{g}}_t \cdot B_{t-1}/B_t$ , we have, for all  $x > 0$ ,

$$\begin{aligned}
 \sum_{i=1}^t -2R\|\hat{\mathbf{g}}_i\| + \left\| \sum_{i=1}^t \hat{\mathbf{g}}_i \right\| x &\leq \sum_{i=1}^t \langle \hat{\mathbf{g}}_i, \mathbf{x}_i - \bar{\mathbf{x}}_{i-1} \rangle + \left\| \sum_{i=1}^t \hat{\mathbf{g}}_i \right\| x, \quad (\text{since } \mathbf{x}_i, \bar{\mathbf{x}}_{i-1} \in \mathcal{C}) \\
 &\leq \sum_{i=1}^t \langle \hat{\mathbf{g}}_i, \mathbf{w}_i - \bar{\mathbf{x}}_{i-1} \rangle + \left\| \sum_{i=1}^t \hat{\mathbf{g}}_i \right\| x + \sum_{i=1}^t \xi_i(\mathbf{x}_i) + (2\delta_i + \Delta_i)R\|\mathbf{g}_i\|, \\
 &= \sum_{i=1}^t \langle \hat{\mathbf{g}}_i, \mathbf{u}_i \rangle + \left\| \sum_{i=1}^t \hat{\mathbf{g}}_i \right\| x + \sum_{i=1}^t (2\delta_i + \Delta_i)R\|\mathbf{g}_i\| + \sum_{i=1}^t \xi_i(\mathbf{x}_i), \\
 &\leq 2x \sqrt{\left( \epsilon^2 + \sum_{i=1}^t \|\hat{\mathbf{g}}_i\|^2 \right) \ln_+ \left( \frac{2\zeta_t^2 x Q_t}{\epsilon^2} \right)} + \sum_{i=1}^t \xi_i(\mathbf{x}_i) \\
 &\quad + \sum_{i=1}^t (2\delta_i + \Delta_i)R\|\mathbf{g}_i\| + 4\zeta'_t x \ln \left( \frac{4\zeta_t^2 B_T x \sqrt{Q_t}}{\epsilon^2} \right) + \epsilon, \quad (118)
 \end{aligned}$$

where  $\mathbf{u}_i$  is the  $i$ th output of FreeGrad,  $Q_t := \epsilon^2 + \sum_{i=1}^t \|\mathbf{g}_i\|^2 \leq \epsilon^2 + tB_t^2$ , and  $\zeta'_t := \max_{s \in [t]} \{(1 + d\kappa + d\Delta_s)B_s\}$ . The passage to (118) follows from the regret bound FreeGrad (see Proposition 26) and the fact that  $\|\hat{\mathbf{g}}_t\| \leq \zeta_t \|\mathbf{g}_t\|$  (see Lemma 34). Moving forward, we let  $G_t := \epsilon + \sum_{i=1}^t (2\delta_i + \Delta_i)R\|\mathbf{g}_i\|$ . Dividing by  $x$  in (118) and solving for  $\left\| \sum_{i=1}^t \hat{\mathbf{g}}_i \right\|$ , we get

$$\begin{aligned}
 \left\| \sum_{i=1}^t \hat{\mathbf{g}}_i \right\| &\leq 2 \sqrt{\left( \epsilon^2 + \sum_{i=1}^t \|\hat{\mathbf{g}}_i\|^2 \right) \ln_+ \left( \frac{2\zeta_t^2 x Q_t}{\epsilon^2} \right)} \\
 &\quad + 4\zeta'_t \ln \left( \frac{4\zeta_t^2 B_t x \sqrt{Q_t}}{\epsilon^2} \right) + \frac{G_t}{x} + \frac{2R}{x} \sum_{i=1}^t \|\hat{\mathbf{g}}_i\| + \frac{1}{x} \sum_{i=1}^t \xi_i(\mathbf{x}_i).
 \end{aligned}$$

Set  $x = R + 2R \sum_{i=1}^t (\|\hat{\mathbf{g}}_i\| \vee \|\mathbf{g}_i\|) / B_T$  and using the facts that  $\|\hat{\mathbf{g}}_t\| \leq \zeta_t \|\mathbf{g}_t\|$  and  $|\xi_t(\mathbf{x})| \leq 4d(\kappa + \delta_t)R\|\mathbf{g}_t\|$ , for all  $\mathbf{x} \in \mathcal{C}$  (see Lemma 34), we conclude that:

$$\begin{aligned}
 \left\| \sum_{i=1}^t \hat{\mathbf{g}}_i \right\| &\leq M_t \sqrt{\epsilon^2 + \sum_{i=1}^t \|\hat{\mathbf{g}}_i\|^2} + N_t, \quad \text{where} \\
 M_t &:= 2 \sqrt{\ln_+ \left( \frac{2\zeta_t^2 (R + 2\zeta_t R t) Q_t}{\epsilon^2} \right)}, \\
 \text{and } N_t &:= 4\zeta'_t \cdot \ln_+ \left( \frac{4\zeta_t^2 B_t (R + 2\zeta_t R t) \sqrt{Q_t}}{\epsilon^2} \right) + B_t + \frac{G_t}{R} + 2d(\kappa + \delta)B_t.
 \end{aligned}$$

We recall that  $Q_t := \epsilon^2 + \sum_{i=1}^t \|\mathbf{g}_i\|^2 \leq \epsilon^2 + tB_t^2$ . With this in hand, we have

$$\begin{aligned}
 \sum_{t=1}^T \langle \hat{\mathbf{g}}_t, \bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_T \rangle &\leq \sum_{t=1}^T \left\| \sum_{i=1}^t \hat{\mathbf{g}}_i \right\| \|\bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_T\|, \\
 &\leq M_T \sum_{t=1}^T \sqrt{\epsilon^2 + \sum_{i=1}^t \|\hat{\mathbf{g}}_i\|^2} \|\bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_T\| + N_T \sum_{t=1}^T \|\bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_T\|.
 \end{aligned}$$

Now, we relate  $\|\bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_t\|$  to  $\|\mathbf{x}_t - \bar{\mathbf{x}}_t\|$ :

$$\bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_t = \bar{\mathbf{x}}_{t-1} - \frac{Z_{t-1}\bar{\mathbf{x}}_{t-1} + \|\hat{\mathbf{g}}_t\|^2 \mathbf{x}_t}{Z_t} = \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_t} (\bar{\mathbf{x}}_{t-1} - \mathbf{x}_t) = \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_t} (\bar{\mathbf{x}}_t - \mathbf{x}_t) + \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_t} (\bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_t),$$

that implies

$$Z_t \cdot (\bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_t) = \|\hat{\mathbf{g}}_t\|^2 (\mathbf{x}_t - \bar{\mathbf{x}}_t) + \|\hat{\mathbf{g}}_t\|^2 (\bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_t),$$

that is

$$\bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_t = \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_{t-1}} (\mathbf{x}_t - \bar{\mathbf{x}}_t). \quad (119)$$

Hence, we have

$$M_T \sum_{t=1}^T \sqrt{\epsilon^2 + \sum_{i=1}^t \|\hat{\mathbf{g}}_i\|^2} \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t-1}\| \leq M_T \sum_{t=1}^T \sqrt{Z_t} \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_{t-1}} \|\mathbf{x}_t - \bar{\mathbf{x}}_t\|,$$

and

$$N_T \sum_{t=1}^T \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t-1}\| \leq N_T \sum_{t=1}^T \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_{t-1}} \|\mathbf{x}_t - \bar{\mathbf{x}}_t\| \leq 2RN_T \sum_{t=1}^T \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_{t-1}}.$$

Using Cauchy–Schwarz inequality, we have

$$M_T \sum_{t=1}^T \sqrt{Z_t} \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_{t-1}} \|\mathbf{x}_t - \bar{\mathbf{x}}_t\| \leq M_T \sqrt{\sum_{t=1}^T \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_{t-1}}} \sqrt{\sum_{t=1}^T \frac{Z_t}{Z_{t-1}} \|\hat{\mathbf{g}}_t\|^2 \|\mathbf{x}_t - \bar{\mathbf{x}}_t\|^2}.$$

So, putting together the last inequalities, we have

$$\sum_{t=1}^T \langle \hat{\mathbf{g}}_t, \bar{\mathbf{x}}_{t-1} - \bar{\mathbf{x}}_T \rangle \leq M_T \sqrt{\sum_{t=1}^T \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_{t-1}}} \sqrt{\sum_{t=1}^T \frac{Z_t}{Z_{t-1}} \|\hat{\mathbf{g}}_t\|^2 \|\mathbf{x}_t - \bar{\mathbf{x}}_t\|^2} + 2RN_T \sum_{t=1}^T \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_{t-1}}.$$

We now focus on the the term  $\sum_{t=1}^T \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_{t-1}}$  that is easily bounded:

$$\sum_{t=1}^T \frac{\|\hat{\mathbf{g}}_t\|^2}{Z_{t-1}} \leq \sum_{t=1}^T \frac{\|\hat{\mathbf{g}}_t\|^2}{\epsilon^2 + \sum_{i=1}^t \|\hat{\mathbf{g}}_i\|^2} \leq \ln \frac{Z_T}{\epsilon^2},$$

where the first inequality follows by the fact that  $\tilde{B}_{t-1} \geq \|\hat{\mathbf{g}}_t\|$ , and in the last inequality we used the inequality

$$\sum_{t=1}^T \frac{a_t}{\sum_{i=0}^t a_i} \leq \ln \left( \frac{\sum_{t=0}^T a_t}{a_0} \right),$$

for all  $a_t \geq 0$ . To bound the term  $\sum_{t=1}^T \frac{Z_t}{Z_{t-1}} \|\hat{\mathbf{g}}_t\|^2 \|\mathbf{x}_t - \bar{\mathbf{x}}_t\|^2$  from above, observe that

$$\begin{aligned}
 \sigma_T^2 Z_T &= (\epsilon^2 + \tilde{B}_T^2) \|\bar{\mathbf{x}}_T\|^2 + \sum_{t=1}^T \|\hat{\mathbf{g}}_t\|^2 \|\mathbf{x}_t - \bar{\mathbf{x}}_t\|^2, \\
 &\geq (\epsilon^2 + \tilde{B}_{T-1}^2) \|\bar{\mathbf{x}}_T\|^2 + \sum_{t=1}^{T-1} \|\hat{\mathbf{g}}_t\|^2 \|\mathbf{x}_t - \bar{\mathbf{x}}_t\|^2 + \|\hat{\mathbf{g}}_T\|^2 \|\mathbf{x}_T - \bar{\mathbf{x}}_T\|^2, \\
 &= Z_{T-1} \cdot (\sigma_{T-1}^2 + \|\bar{\mathbf{x}}_T - \bar{\mathbf{x}}_{T-1}\|^2) + \|\hat{\mathbf{g}}_T\|^2 \|\mathbf{x}_T - \bar{\mathbf{x}}_T\|^2, \\
 &= Z_{T-1} \sigma_{T-1}^2 + \|\hat{\mathbf{g}}_T\|^2 \left(1 + \frac{\|\hat{\mathbf{g}}_T\|^2}{Z_{T-1}}\right) \|\mathbf{x}_T - \bar{\mathbf{x}}_T\|^2, \\
 &= Z_{T-1} \sigma_{T-1}^2 + \|\hat{\mathbf{g}}_T\|^2 \frac{Z_T}{Z_{T-1}} \|\mathbf{x}_T - \bar{\mathbf{x}}_T\|^2,
 \end{aligned}$$

where the third equality comes from bias-variance decomposition and the fourth one comes from (119). Hence, we have

$$\sum_{t=1}^T \frac{Z_t}{Z_{t-1}} \|\hat{\mathbf{g}}_t\|^2 \|\mathbf{x}_t - \bar{\mathbf{x}}_t\|^2 = \sum_{t=1}^T (\sigma_t^2 Z_t - \sigma_{t-1}^2 Z_{t-1}) \leq \sigma_T^2 Z_T.$$

Putting all together, we have the stated bound.  $\blacksquare$

## I.6. Proof of Theorem 18 (The Smooth Stochastic Case)

We will need to use the result of Lemma 31 (in particular (73)) to show the desired convergence rate. In order for the result of Lemma 31 to be valid in the setting of Theorem 18, we need to show that the iterates  $(\mathbf{w}_t)$  in Algorithm 6 are in  $\mathcal{B}(R)$ , which is what we do next:

**Lemma 37** *In the setting of Theorem 18, we have  $(\mathbf{w}_t) \subset \mathcal{B}(R)$ .*

**Proof** Let  $\gamma_t$  be as in Algorithm 6. For  $t = 1$ , we have  $\mathbf{w}_1 = \mathbf{0}$ . For  $t > 1$ , we have  $\mathbf{w}_t = (1 - \mu_t)(\mathbf{x}_{t-1} - \eta_{t-1} \mathbf{g}_{t-1}) + \mu_t \mathbf{u}_t$  and  $\|\mathbf{u}_t\| \leq R'$  (since it is the output of FTRL-prox with parameter  $R'$ ), and so

$$\begin{aligned}
 \|\mathbf{w}_t\| &\leq (1 - \mu_t) \|\mathbf{x}_{t-1}\| + \mu_t R' + |\eta_{t-1}| \|\mathbf{g}_{t-1}\| = (1 - \mu_t) \|\mathbf{x}_{t-1}\| + \mu_t R' + \frac{\nu R' \|\mathbf{g}_{t-1}\|}{\sqrt{Z_{t-1}}}, \\
 &\leq (1 - \mu_t) \|\mathbf{x}_{t-1}\| + (\mu_t + \nu) R'.
 \end{aligned}$$

Thus,  $\mathbf{w}_t \in \mathcal{B}(R)$  if  $\mathbf{x}_t \in \mathcal{C}$  (which implies  $\|\mathbf{x}_t\| \leq R'$ ). We will now show that  $\mathbf{x}_t \in \mathcal{C}$ . We consider two cases. Suppose that  $\gamma_t < 1$ . In this case,  $\mathbf{x}_t = \mathbf{w}_t$ . If  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \leq 9/16$ , then by definition of the Gauge function we have  $\mathbf{w}_t \in \mathcal{C}$  and the same holds for  $\mathbf{x}_t$  (since  $\mathbf{x}_t = \mathbf{w}_t$ ). On the other hand, if  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \geq 9/16$ , then by Lemma 10, we have  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \leq \gamma_t < 1$ , and so  $\mathbf{x}_t = \mathbf{w}_t \in \mathcal{C}$ .

Now suppose that  $\gamma_t \geq 1$ . In this case, we have  $\mathbf{x}_t = \mathbf{w}_t / \gamma_t$  and so  $\gamma_{\mathcal{C}}(\mathbf{x}_t) = \gamma_{\mathcal{C}}(\mathbf{w}_t / \gamma_t) = \gamma_{\mathcal{C}}(\mathbf{w}_t) / \gamma_t \leq \gamma_{\mathcal{C}}(\mathbf{w}_t) / \gamma_{\mathcal{C}}(\mathbf{w}_t) = 1$  (since  $\gamma_{\mathcal{C}}(\mathbf{w}_t) \leq \gamma_t$  by Lemma 10), and so  $\mathbf{x}_t \in \mathcal{C}$ .  $\blacksquare$

The following lemma will also be useful to us in the proof of Theorem 18:



**Lemma 38** Let  $\beta > 0$ ,  $t \geq 1$ ,  $\boldsymbol{\xi}_t \in \mathbb{R}^d$  be a random vector satisfying (23) for  $\sigma > 0$ , and  $\ell_t$  be as in (22). Further, suppose that  $\mathcal{C}$  satisfies (26). When  $f$  is  $\beta$ -smooth on  $\mathcal{C}$ , we have,

$$\mathbb{E}[\|\mathbf{g}_t\|^2] \leq \sigma^2 + 2R'\beta, \quad \text{for any } \mathbf{x} \in \mathcal{C} \text{ and } \mathbf{g}_t \in \partial\ell_t(\mathbf{x}).$$

**Proof** Let  $\mathbf{x} \in \mathcal{C}$  and  $\mathbf{g}_t \in \partial\ell_t(\mathbf{x})$ . First, since  $f$  is differentiable,  $\ell_t$  is also differentiable. Thus,  $\mathbf{g}_t = \nabla f(\mathbf{x}) + \boldsymbol{\xi}_t$ , and so

$$\mathbb{E}[\|\mathbf{g}_t\|^2] = \mathbb{E}[\|\nabla f(\mathbf{x})\|^2 + 2\langle \nabla f(\mathbf{x}), \boldsymbol{\xi}_t \rangle + \|\boldsymbol{\xi}_t\|^2] = \|\nabla f(\mathbf{x})\|^2 + \mathbb{E}[\|\boldsymbol{\xi}_t\|^2] \leq 2R'\beta + \sigma^2,$$

where the last inequality follows (23) and (25).  $\blacksquare$

We now present the proof of Theorem 18. We note that the proof is very similar to that of (Cutkosky, 2019, Theorem 4) with modifications to account for the application of our projection-free reduction (Algorithm 1).

**Proof of Theorem 18.** Throughout this proof, we let  $\lambda_t = t$ ,  $\forall t$ , and  $\nu$  and  $R'$  be as in (26). Note that  $\Lambda_t$  in Algorithm 6 satisfies  $\Lambda_t = \sum_{s=1}^t \lambda_s$ . By a standard convexity argument, we have

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^T \lambda_t (f(\mathbf{x}_t) - f(\mathbf{x})) \right] &\leq \mathbb{E} \left[ \sum_{t=1}^T \lambda_t \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x} \rangle \right] \\ &\leq \mathbb{E} \left[ \sum_{t=1}^T \lambda_t \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{x} \rangle \right] + \mathbb{E} \left[ \sum_{t=1}^T 3\lambda_t \delta_t R \|\mathbf{g}_t\| \right], \quad (\text{by Lemma 31}) \\ &= \mathbb{E} \left[ \sum_{t=1}^T \lambda_t \langle \tilde{\mathbf{g}}_t, \mathbf{u}_t - \mathbf{x} \rangle + \sum_{t=1}^T \lambda_t \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{u}_t \rangle \right] \\ &\quad + \mathbb{E} \left[ \sum_{t=1}^T 3\lambda_t \delta_t R \cdot \mathbb{E}[\|\mathbf{g}_t\| \mid \mathbf{x}_t] \right], \\ &\leq \mathbb{E} \left[ \sum_{t=1}^T \lambda_t \langle \tilde{\mathbf{g}}_t, \mathbf{u}_t - \mathbf{x} \rangle + \sum_{t=1}^T \lambda_t \langle \tilde{\mathbf{g}}_t, \mathbf{w}_t - \mathbf{u}_t \rangle \right] \\ &\quad + 3R\sqrt{\sigma^2 + 2\beta R'} \sum_{t=1}^T \lambda_t \delta_t, \end{aligned} \tag{120}$$

where the last inequality follows by Lemma 38 and Jensen's inequality. By letting  $\mathbf{y}_s := \mathbf{x}_s - \eta_s \mathbf{g}_s$ , we have by Line 8 of Algorithm 6

$$\lambda_s \cdot (\mathbf{w}_s - \mathbf{u}_s) = \Lambda_{s-1} \cdot (\mathbf{y}_{s-1} - \mathbf{w}_s), \quad \text{for all } s \geq 1. \tag{121}$$

Moreover, the first term on the RHS of (120) is the regret of the FTRL-prox instance (Algorithm A) within Algorithm 6. Thus, by Proposition 25 and Lemma 31, this term is bounded from above by

$$\mathcal{R}_T^{\text{FTRL}}(\mathbf{x}) := 4(1 + \kappa)R' \sqrt{\sum_{t=1}^T \lambda_t^2 \|\mathbf{g}_t\|^2} + 4R' \sum_{t=1}^T \lambda_t \Delta_t \|\mathbf{g}_t\|, \quad \forall \mathbf{x} \in \mathcal{C}. \tag{122}$$

Plugging (121) and (122) into (120), yields

$$\begin{aligned}
 & \mathbb{E} \left[ \sum_{t=1}^T \lambda_t (f(\mathbf{x}_t) - f(\mathbf{x})) \right] \\
 & \leq \mathbb{E} \left[ \mathcal{R}_T^{\text{FTRL}}(\mathbf{x}) \right] + \mathbb{E} \left[ \sum_{t=1}^T \Lambda_{t-1} \langle \tilde{\mathbf{g}}_t, \mathbf{y}_{t-1} - \mathbf{x}_t \rangle \right] + 3R\sqrt{\sigma^2 + 2\beta R'} \sum_{t=1}^T \lambda_t \delta_t, \\
 & \leq \mathbb{E} \left[ \mathcal{R}_T^{\text{FTRL}}(\mathbf{x}) \right] + \mathbb{E} \left[ \sum_{t=1}^T \Lambda_{t-1} \langle \mathbf{g}_t, \mathbf{y}_{t-1} - \mathbf{x}_t \rangle \right] \\
 & \quad + \mathbb{E} \left[ \sum_{t=1}^T 3\Lambda_{t-1} \delta_t R \|\mathbf{g}_t\| \right] + 3R\sqrt{\sigma^2 + 2\beta R'} \sum_{t=1}^T \lambda_t \delta_t, \tag{123} \\
 & = \mathbb{E} \left[ \mathcal{R}_T^{\text{FTRL}}(\mathbf{x}) \right] + \mathbb{E} \left[ \sum_{t=1}^T \Lambda_{t-1} \langle \mathbf{g}_t, \mathbf{y}_{t-1} - \mathbf{x}_t \rangle \right] \\
 & \quad + \mathbb{E} \left[ \sum_{t=1}^T 3\delta_t \Lambda_{t-1} R \cdot \mathbb{E}[\|\mathbf{g}_t\| \mid \mathbf{x}_t] \right] + 3R\sqrt{\sigma^2 + 2\beta R'} \sum_{t=1}^T \lambda_t \delta_t, \\
 & \leq \mathbb{E} \left[ \mathcal{R}_T^{\text{FTRL}}(\mathbf{x}) \right] + \mathbb{E} \left[ \sum_{t=1}^T \Lambda_{t-1} \langle \mathbf{g}_t, \mathbf{y}_{t-1} - \mathbf{x}_t \rangle \right] + 3R\sqrt{\sigma^2 + 2\beta R'} \sum_{t=1}^T (\lambda_t + \Lambda_{t-1}) \delta_t,
 \end{aligned}$$

where (123) follows by Lemma 31, and the last inequality follows by Lemma 38 and Jensen's inequality. Next, we use convexity again to argue

$$\mathbb{E}[\langle \mathbf{g}_t, \mathbf{y}_{t-1} - \mathbf{x}_t \rangle] \leq \mathbb{E}[f(\mathbf{y}_{t-1}) - f(\mathbf{x}_t)],$$

and then we subtract  $\mathbb{E}[\sum_{t=1}^T \lambda_t f(\mathbf{x}_t)]$  from both sides:

$$\mathbb{E}[-\Lambda_T f(\mathbf{x})] \leq \mathbb{E} \left[ \mathcal{R}_T^{\text{FTRL}}(\mathbf{x}) \right] + \sum_{t=1}^T (\Lambda_{t-1} f(\mathbf{y}_{t-1}) - \Lambda_t f(\mathbf{x}_t)) + 3R\sqrt{\sigma^2 + 2\beta R'} \sum_{t=1}^T (\lambda_t + \Lambda_{t-1}) \delta_t.$$

Now we use smoothness to relate  $f(\mathbf{y}_t)$  to  $f(\mathbf{x}_t)$ . Let  $\nu := 4\sqrt{2}(1 + \kappa)$  so that  $\eta_t = \nu R' / \sqrt{Z_t}$ , where  $Z_t$  is as in Algorithm 6; i.e.  $Z_t := \varepsilon^2 + \sum_{s=1}^t \Lambda_s \|\mathbf{g}_s\|^2$ . With this, we have:

$$\begin{aligned}
 \mathbb{E}[f(\mathbf{y}_t)] & \leq \mathbb{E}[f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)(\mathbf{y}_t - \mathbf{x}_t) + \frac{\beta}{2} \|\mathbf{x}_t - \mathbf{y}_t\|^2], \\
 & \leq \mathbb{E} \left[ f(\mathbf{x}_t) - \eta_t \|\mathbf{g}_t\|^2 + \eta_t \langle \boldsymbol{\xi}_t, \mathbf{g}_t \rangle + \frac{\beta \eta_t^2 \|\mathbf{g}_t\|^2}{2} \right].
 \end{aligned}$$

Then multiply by  $\Lambda_t$ :

$$\mathbb{E}[\Lambda_t (f(\mathbf{y}_t) - f(\mathbf{x}_t))] \leq \mathbb{E} \left[ -\frac{\nu R' \Lambda_t \|\mathbf{g}_t\|^2}{\sqrt{\varepsilon^2 + \sum_{i=1}^t \Lambda_i \|\mathbf{g}_i\|^2}} + \frac{\beta \eta_t^2 \Lambda_t \|\mathbf{g}_t\|^2}{2} + \eta_t \Lambda_t \langle \boldsymbol{\xi}_t, \mathbf{g}_t \rangle \right].$$

Next, we make use of the following facts (see e.g. (Cutkosky, 2019; Levy et al., 2018)): for positive numbers  $\alpha_0, \dots, \alpha_n$ ,

$$\sqrt{\sum_{i=1}^n \alpha_i} \leq \sum_{i=1}^n \frac{\alpha_i}{\sqrt{\sum_{j=1}^i \alpha_j}} \leq 2 \sqrt{\sum_{i=1}^n \alpha_i} \quad \text{and} \quad \sum_{i=1}^n \frac{\alpha_i}{\alpha_0 + \sum_{j=1}^i \alpha_j} \leq \ln \left( \alpha_0 + \sum_{i=1}^n \alpha_i \right) - \ln \alpha_0.$$

Using this, we obtain

$$\mathbb{E} \left[ \sum_{t=1}^T \Lambda_t (f(\mathbf{y}_t) - f(\mathbf{x}_t)) \right] \tag{124}$$

$$\begin{aligned} &\leq \mathbb{E} \left[ -\nu R' \sqrt{\varepsilon^2 + \sum_{t=1}^T \Lambda_t \|\mathbf{g}_t\|^2} + \frac{\nu^2 (R')^2 \beta}{2} \ln \left( 1 + \frac{\sum_{t=1}^T \Lambda_t \|\mathbf{g}_t\|^2}{\varepsilon^2} \right) \right] \\ &\quad + \nu R' \varepsilon + \mathbb{E} \left[ \sum_{t=1}^T \eta_t \langle \boldsymbol{\xi}_t, \Lambda_t \mathbf{g}_t \rangle \right], \\ &\leq \mathbb{E} \left[ -\nu R' \sqrt{\varepsilon^2 + \sum_{t=1}^T \Lambda_t \|\mathbf{g}_t\|^2} \right] + \nu R' \varepsilon + \mathbb{E} \left[ \sum_{t=1}^T \eta_t \langle \boldsymbol{\xi}_t, \Lambda_t \mathbf{g}_t \rangle \right] \\ &\quad + \frac{\nu^2 (R')^2 \beta}{2} \ln \left( 1 + \frac{\sum_{t=1}^T \Lambda_t \mathbb{E} [\|\mathbf{g}_t\|^2]}{\varepsilon^2} \right), \end{aligned} \tag{125}$$

$$\begin{aligned} &\leq \mathbb{E} \left[ -\nu R' \sqrt{\varepsilon^2 + \sum_{t=1}^T \Lambda_t \|\mathbf{g}_t\|^2} \right] + \nu R' \varepsilon + \frac{\nu^2 (R')^2 \beta}{2} \ln \left( 1 + \frac{(\sigma^2 + 2R'\beta) \sum_{t=1}^T \Lambda_t}{\varepsilon^2} \right) \\ &\quad + \mathbb{E} \left[ \sum_{t=1}^T \eta_t \langle \boldsymbol{\xi}_t, \Lambda_t \mathbf{g}_t \rangle \right], \end{aligned} \tag{126}$$

where (125) follows by Jensen's inequality (the log is concave) and the triangular inequality, and (126) follows by Lemma 38. Using Cauchy-Schwarz, we obtain:

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^T \eta_t \langle \boldsymbol{\xi}_t, \Lambda_t \mathbf{g}_t \rangle \right] &\leq \mathbb{E} \left[ \sqrt{\sum_{t=1}^T \Lambda_t \|\boldsymbol{\xi}_t\|^2} \sqrt{\sum_{t=1}^T \eta_t^2 \Lambda_t \|\mathbf{g}_t\|^2} \right], \\ &\leq \mathbb{E} \left[ \nu R' \sqrt{\sum_{t=1}^T \Lambda_t \|\boldsymbol{\xi}_t\|^2} \cdot \ln \left( 1 + \varepsilon^{-2} \sum_{s=1}^T \Lambda_s \|\mathbf{g}_s\|^2 \right) \right], \\ &\leq \mathbb{E} \left[ \nu R' \sqrt{\sum_{t=1}^T \Lambda_t \|\boldsymbol{\xi}_t\|^2} \cdot \ln \left( 1 + 2\varepsilon^{-2} \sum_{s=1}^T \Lambda_s (\|\boldsymbol{\xi}_s\|^2 + \|\nabla f(\mathbf{x}_s)\|^2) \right) \right], \\ &\leq \mathbb{E} \left[ \nu R' \sqrt{\sum_{t=1}^T \Lambda_t \|\boldsymbol{\xi}_t\|^2} \cdot \ln \left( 1 + 2\varepsilon^{-2} \sum_{s=1}^T \Lambda_s (\|\boldsymbol{\xi}_s\|^2 + 2R'\beta) \right) \right], \quad (\text{by (25)}), \end{aligned}$$

$$\begin{aligned}
 &\leq \mathbb{E} \left[ \nu R' \sqrt{\sum_{t=1}^T \Lambda_t \|\boldsymbol{\xi}_t\|^2 \cdot \ln \left( 1 + 2\varepsilon^{-2} \sum_{s=1}^T \Lambda_s (\|\boldsymbol{\xi}_s\|^2 + 2R'\beta) \right)} \right], \quad (127) \\
 &\leq \nu \sigma R' \sqrt{\sum_{t=1}^T \Lambda_t \cdot \ln \left( 1 + \frac{2(\sigma^2 + 2R'\beta) \sum_{s=1}^T \Lambda_s}{\varepsilon^2} \right)},
 \end{aligned}$$

where (127) follows by the concavity of the function  $x \mapsto \sqrt{x \cdot \ln(a + bx)}$ , for any  $a \geq 1, b > 0$ , and  $x > 0$  (see Lemma 41) and Jensen's inequality, and the last inequality follows by the fact that  $\mathbb{E}[\|\boldsymbol{\xi}_t\|^2] \leq \sigma^2$ .

Combining everything, we get

$$\begin{aligned}
 \mathbb{E} \left[ \sum_{t=1}^T -\lambda_t f(\mathbf{x}) \right] &\leq \mathbb{E} \left[ \mathcal{R}_T^{\text{FTRL}}(\mathbf{x}) + \sum_{t=1}^T (\Lambda_{t-1} f(\mathbf{y}_{t-1}) - \Lambda_t f(\mathbf{y}_t)) \right] \\
 &\quad + \frac{\nu^2 \beta (R')^2}{2} \ln \left( 1 + \frac{\sum_{t=1}^T \Lambda_t (\sigma^2 + 2R'\beta)}{\varepsilon^2} \right) - \nu R' \sqrt{\varepsilon^2 + \sum_{t=1}^T \Lambda_t \|\mathbf{g}_t\|^2} \\
 &\quad + \nu R' \varepsilon + \nu R' \sigma \sqrt{\sum_{t=1}^T \Lambda_t \ln \left( 1 + \frac{2(\sigma^2 + 2R'\beta) \sum_{s=1}^T \Lambda_s}{\varepsilon^2} \right)} \\
 &\quad + 3R \sqrt{\sigma^2 + 2\beta R'} \sum_{t=1}^T (\lambda_t + \Lambda_{t-1}) \delta_t.
 \end{aligned}$$

Now observe that  $t^2 > \Lambda_t > \lambda_t^2/2$  and recall  $\mathcal{R}_T^{\text{FTRL}}(\mathbf{x}) = 4(1 + \kappa)R' \sqrt{\sum_{t=1}^T \lambda_t^2 \|\mathbf{g}_t\|^2} + 4R' \sum_{t=1}^T \lambda_t \Delta_t \|\mathbf{g}_t\|$ . Therefore, since  $\nu = 4\sqrt{2}(1 + \kappa)$  we have:

$$\begin{aligned}
 &\mathbb{E} \left[ \mathcal{R}_T^{\text{FTRL}}(\mathbf{x}) - \nu R' \sqrt{\varepsilon^2 + \sum_{t=1}^T \Lambda_t \|\mathbf{g}_t\|^2} \right] \\
 &\leq \mathbb{E} \left[ 4(1 + \kappa)R' \sqrt{\sum_{t=1}^T \lambda_t^2 \|\mathbf{g}_t\|^2} - 4\sqrt{2}(1 + \kappa)R' \sqrt{\sum_{t=1}^T \lambda_t^2 \|\mathbf{g}_t\|^2/2} \right] + \mathbb{E} \left[ 4R' \sum_{t=1}^T \lambda_t \Delta_t \|\mathbf{g}_t\| \right], \\
 &= \mathbb{E} \left[ 4R' \sum_{t=1}^T \lambda_t \Delta_t \cdot \mathbb{E}[\|\mathbf{g}_t\| \mid \mathbf{x}_t] \right] \leq 4R' \sqrt{\sigma^2 + 2\beta R'} \sum_{t=1}^T \lambda_t \delta_t,
 \end{aligned}$$

where the last inequality follows by Lemma 38. Also, observe that  $\sum_{t=1}^T \Lambda_t \leq \sum_{t=1}^T t^2 \leq T^3$ . Thus, we telescope the sum to obtain:

$$\begin{aligned}
 &\mathbb{E}[\Lambda_T (f(\mathbf{y}_T) - f(\mathbf{x}))] \\
 &\leq \nu R' \varepsilon + \frac{\nu^2 (R')^2 \beta}{2} \ln \left( 1 + \frac{(\sigma^2 + 2R'\beta)T^3}{\varepsilon^2} \right) + \nu R' T^{3/2} \sigma \sqrt{\ln \left( 1 + \frac{2(\sigma^2 + 2R'\beta)T^3}{\varepsilon^2} \right)} \\
 &\quad + 3R \sqrt{\sigma^2 + 2\beta R'} \sum_{t=1}^T (\Lambda_{t-1} + 3\lambda_t) \delta_t,
 \end{aligned}$$

$$\begin{aligned} &\leq \nu R' \varepsilon + \frac{\nu^2 (R')^2 \beta}{2} \ln \left( 1 + \frac{(\sigma^2 + 2R' \beta) T^3}{\varepsilon^2} \right) + \nu R' T^{3/2} \sigma \sqrt{\ln \left( 1 + \frac{2(\sigma^2 + 2R' \beta) T^3}{\varepsilon^2} \right)} \\ &\quad + 3\delta R (\ln T + 6) \sqrt{\sigma^2 + 2\beta R'}, \end{aligned} \quad (128)$$

where the last inequality follows by the fact that  $\delta_t = \delta/t^3$ ,  $\sum_{t=1}^T 1/t \leq \ln T$ ; and  $\sum_{t=1}^T 1/t^2 \leq 2$ . Dividing (128) by  $\Lambda_T = \sum_{t=1}^T \lambda_t = \frac{T(T+1)}{2} > T^2/2$  shows the inequality of theorem. Finally, the fact that  $(\mathbf{x}_t) \subset \mathcal{C}$  follows by Lemma 31.  $\blacksquare$

## Appendix J. Technical Lemmas

This section contains some technical lemmas we need to prove our results.

**Lemma 39** *Let  $(Y_t) \subset \mathbb{R}_{\geq 0}$  be a sequence of random variable satisfying  $\mathbb{E}[Y_t | \mathcal{G}_{t-1}] \leq \delta_t$ , for all  $t \geq 1$ , for some sequence  $(\delta_t) \subset \mathbb{R}_{\geq 0}$ . Then, for any  $\rho \in (0, 1)$  and  $T \geq 1$ , we have with probability at least  $1 - \rho$ ,*

$$\mathbb{P} \left[ \sum_{t=1}^T Y_t \geq (1 + 1/\rho) \sum_{t=1}^T \delta_t \right] \leq \rho.$$

**Proof** Let  $X_t := \sum_{i=1}^t (Y_i - \bar{\delta}_i)$ , where  $\bar{\delta}_i := \mathbb{E}[Y_i | \mathcal{G}_{i-1}] \leq \delta_i$ . The process  $(X_t)$  is a martingale; that is, for all  $i \geq 1$ , we have, for all  $i < t$ ,

$$\mathbb{E}[X_t | \mathcal{G}_i] = \sum_{s=1}^i (Y_s - \bar{\delta}_s) = X_i.$$

Thus, by Doob's martingale inequality (Durrett, 2019, Theorem 4.4.2), we have, for any  $\rho \in (0, 1)$ , and  $T \geq 1$

$$\begin{aligned} \mathbb{P} \left[ \sum_{t=1}^T Y_t \geq (1 + 1/\rho) \sum_{t=1}^T \delta_t \right] &\leq \mathbb{P} \left[ X_T \geq \sum_{t=1}^T \delta_t / \rho \right] \leq \mathbb{P} \left[ \max_{t \leq T} X_t \geq \sum_{t=1}^T \delta_t / \rho \right] \leq \frac{\rho \mathbb{E}[X_T \vee 0]}{\sum_{t=1}^T \delta_t} \\ &\leq \rho. \end{aligned} \quad \blacksquare$$

**Theorem 40** *Let  $\mathcal{F}_1, \dots, \mathcal{F}_n$  be a filtration, and  $X_1, \dots, X_n$  be real random variables such that  $X_i$  is  $\mathcal{F}_i$ -measurable,  $\mathbb{E}[X_i | \mathcal{F}_{i-1}] = 0$ ,  $|X_i| \leq b$ , and  $\sum_{i=1}^n \mathbb{E}[X_i^2 | \mathcal{F}_{i-1}] \leq V$  for some  $b, V \geq 0$ . Then, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,*

$$\sum_{i=1}^n X_i \leq 2\sqrt{V_n \ln(1/\delta)} + b \ln(1/\delta).$$

**Lemma 41** *For any  $a \geq 1$ ,  $b > 0$ , the map  $f : x \mapsto \sqrt{x \cdot \ln(a + bx)}$  is concave for  $x > 0$ .*

**Proof** Let  $a \geq 1$  and  $b > 0$ . We will show that the second derivative of  $f$  is negative for all  $x > 0$ . We have

$$\begin{aligned} \forall x > 0, \quad f''(x) &= \frac{-(a+bx)^2 \log^2(a+bx) + 2abx \log(a+bx) - b^2 x^2}{4(a+bx)^2 (x \log(a+bx))^{3/2}}, \\ &\leq \frac{-a^2 \log^2(a+bx) + 2abx \log(a+bx) - b^2 x^2}{4(a+bx)^2 (x \log(a+bx))^{3/2}}, \\ &= -\frac{-(a \log(a+bx) - bx)^2}{4(a+bx)^2 (x \log(a+bx))^{3/2}} \leq 0. \end{aligned}$$

■

**Lemma 42** *Let  $\mathbf{w} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$  and  $0 < r \leq R$ . Further, let  $\mathcal{C}$  be a closed convex set such that  $\mathcal{B}(r) \subseteq \mathcal{C} \subseteq \mathcal{B}(R)$ . Then, the following properties hold:*

- (a)  $\sigma_{\mathcal{C}^\circ}(\mathbf{w}) = \gamma_{\mathcal{C}}(\mathbf{w})$  and  $(\mathcal{C}^\circ)^\circ = \mathcal{C}$ .
- (b)  $\sigma_{\mathcal{C}}(\alpha \mathbf{w}) = \alpha \sigma_{\mathcal{C}}(\mathbf{w})$  and  $\partial \sigma_{\mathcal{C}}(\alpha \mathbf{w}) = \partial \sigma_{\mathcal{C}}(\mathbf{w}) = \arg \max_{\mathbf{u} \in \mathcal{C}} \langle \mathbf{u}, \mathbf{w} \rangle$ , for all  $\alpha \geq 0$ .
- (c)  $r \|\mathbf{w}\| \leq \sigma_{\mathcal{C}}(\mathbf{w}) \leq R \|\mathbf{w}\|$ ,  $\|\mathbf{w}\|/R \leq \gamma_{\mathcal{C}}(\mathbf{w}) \leq \|\mathbf{w}\|/r$ , and  $\mathcal{B}(1/R) \subseteq \mathcal{C}^\circ \subseteq \mathcal{B}(1/r)$ .
- (d)  $\langle \mathbf{w}, \mathbf{u} \rangle \leq \sigma_{\mathcal{C}}(\mathbf{w}) \cdot \gamma_{\mathcal{C}}(\mathbf{u})$ , for all  $\mathbf{u} \in \mathbb{R}^d$ . (Cauchy Schwarz)
- (e)  $\sigma_{\mathcal{C}}(\mathbf{w} + \mathbf{u}) \leq \sigma_{\mathcal{C}}(\mathbf{w}) + \sigma_{\mathcal{C}}(\mathbf{u})$ , for all  $\mathbf{u} \in \mathbb{R}^d$ . (Sub-additivity)

**Proof** Points (a), (b), and (e) follow from standard results in convex analysis, see e.g. (Molinaro, 2020, Lemma 2) for point (a) and (Hiriart-Urruty and Lemaréchal, 2004) for points (b) and (e). Point (d) follows from (Friedlander et al., 2014, Equation 2.3 & Proposition 2.3) and Point (a). We now show point (c). The first inequality in Point (c) follows by the fact that

$$R \|\mathbf{w}\| \stackrel{(i)}{\geq} \sup_{\mathbf{u} \in \mathcal{C}} \langle \mathbf{u}, \mathbf{w} \rangle = \sigma_{\mathcal{C}}(\mathbf{w}) \geq \inf_{\mathbf{u} \in \mathcal{C}^\circ} \langle \mathbf{u}, \mathbf{w} \rangle \stackrel{(ii)}{\geq} r \|\mathbf{w}\|,$$

where (i) and (ii) follow by the assumption that  $\mathcal{B}(r) \subseteq \mathcal{C} \subseteq \mathcal{B}(R)$ . We now show that  $\mathcal{B}(1/R) \subseteq \mathcal{C}^\circ \subseteq \mathcal{B}(1/r)$ . For any  $\mathbf{x} \in \mathcal{B}(1/R)$ , we have  $\langle \mathbf{x}, \mathbf{w} \rangle \leq 1$  for all  $\mathbf{w} \in \mathcal{C}$ , since  $\mathcal{C} \subseteq \mathcal{B}(R)$ . By definition of the polar set, this implies that  $\mathcal{B}(1/R) \subseteq \mathcal{C}^\circ$ . Now, let  $\mathbf{x} \in \mathcal{C}^\circ$ . This implies that  $\langle \mathbf{x}, \mathbf{w} \rangle \leq 1$  for all  $\mathbf{w} \in \mathcal{C}$ . For  $\mathbf{w} = r\mathbf{x}/\|\mathbf{x}\|$  (which is guaranteed to be in  $\mathcal{C}$  since  $\mathcal{B}(r) \subseteq \mathcal{C}$ ) this inequality implies that  $\|\mathbf{x}\| \leq 1/r$ , and so  $\mathcal{C}^\circ \subseteq \mathcal{B}(1/r)$ . Finally,  $\|\mathbf{w}\|/R \leq \gamma_{\mathcal{C}}(\mathbf{w}) \leq \|\mathbf{w}\|/r$  follows by point (a) and the facts that  $\mathcal{B}(1/R) \subseteq \mathcal{C}^\circ \subseteq \mathcal{B}(1/r)$  and that  $(\mathcal{B}(r) \subseteq \mathcal{C} \subseteq \mathcal{B}(R)) \implies r \|\mathbf{w}\| \leq \sigma_{\mathcal{C}}(\mathbf{w}) \leq R \|\mathbf{w}\|$ , which we just showed. ■